

Deep learning in NLP

Part 3: Transformers, Bert & GPT

Zeinab Rahimi, Feb 2022

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

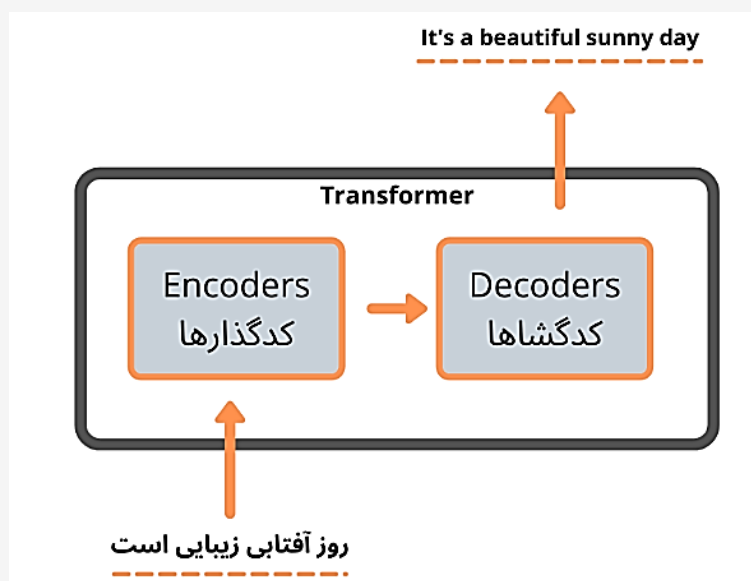
Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez*[†]
University of Toronto
aidan@cs.toronto.edu

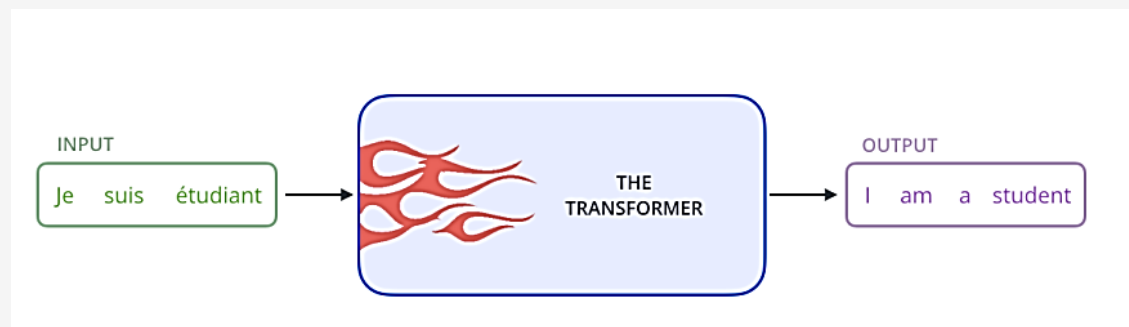
Łukasz Kaiser*
Google Brain
lukaszkaizer@google.com

Transformers

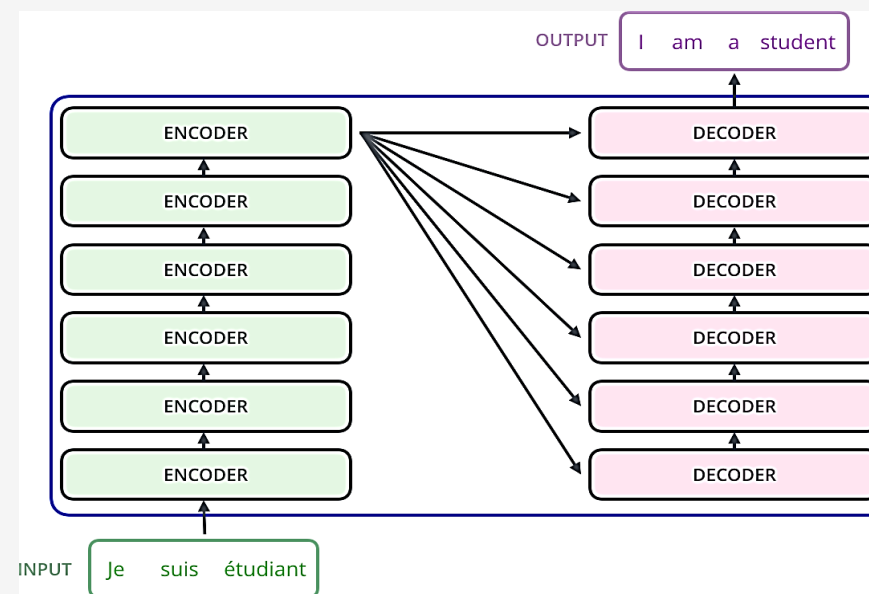
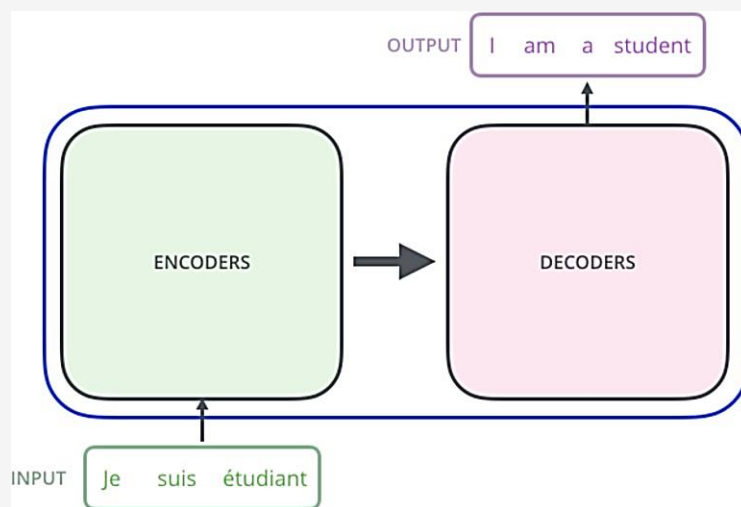


- استفاده از مکانیزم توجه و یادگیری رابطه بین کلمات
- در ساده‌ترین شکل، ترنسفورمر شامل دو مکانیزم جداگانه:
 - یک کدگذار (Encoder) برای خواندن متن ورودی
 - یک کدگشا (Decoder) برای بیان پیش‌بینی محتمل را برای تسک مشخص شده
- لایهٔ انکودر ترنسفورمرها دنباله‌ای از کلمات ورودی را به صورت یکجا می‌خواند.
- برعکس مدل‌های قبلی (RNN و LSTM) که متن ورودی را به ترتیب از چپ به راست یا از راست به چپ می‌خواندند.

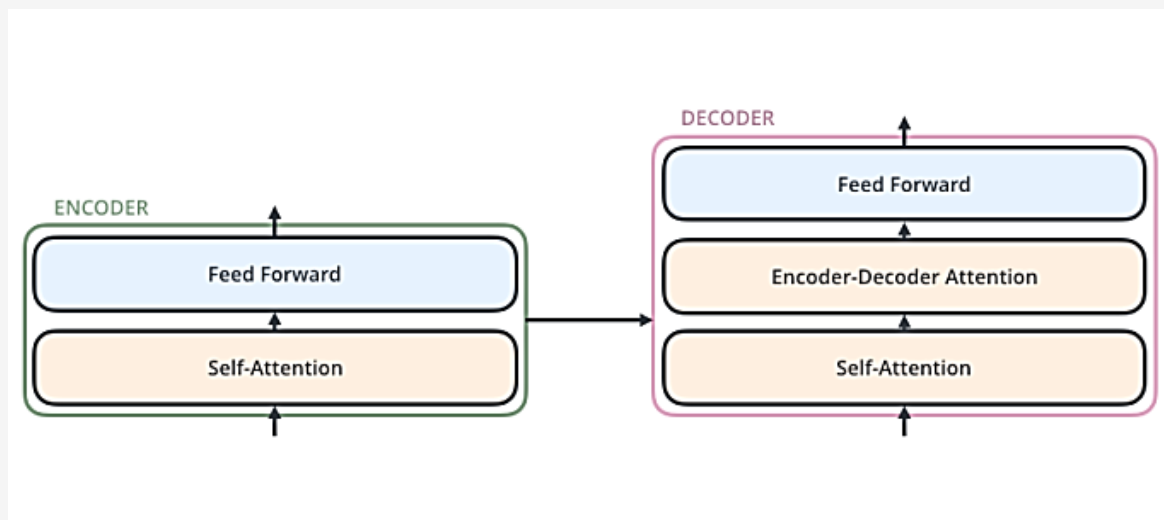
Transformers



- ناکارآمدی زیربنایی شبکه‌های بازگشتی
- در دنباله‌های طولانی و موازی سازی
- معرفی توسط گوگل در سال ۲۰۱۷
- یک بلاک ترنسفورمر: مجموعه‌ای از انکودرها و دیکودرها
- بررسی با یک تسک ترجمه ماشینی
- تعداد لایه‌ها: هاپیرپارامتر



Transformers



• انکودر

- متن ورودی را پردازش می‌کند، قسمت‌های مهم آن را جست‌وجو می‌کند و بر اساس میزان ارتباط هر کلمه با سایر کلمات تشکیل‌دهنده جمله، برای تک‌تک آن‌ها یک تعبیه ایجاد می‌کند.

- یک لایه Self-Attention

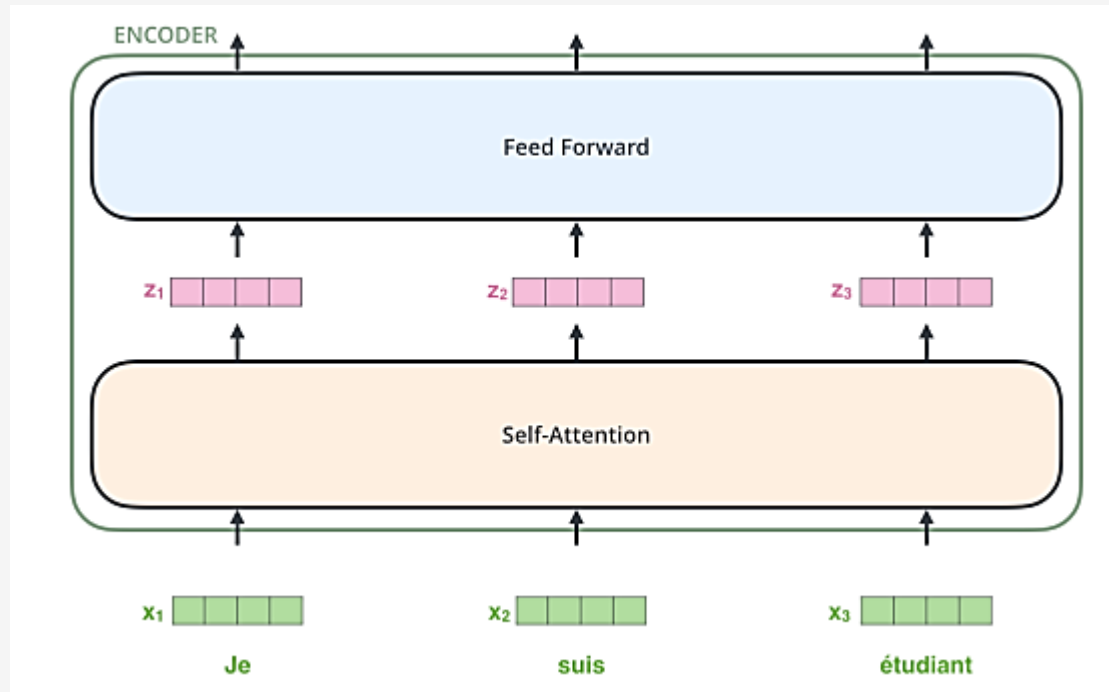
- یک لایه Feed Forward

• دیکودر

- مشابه انکودر + یک لایه Encoder-decoder attention

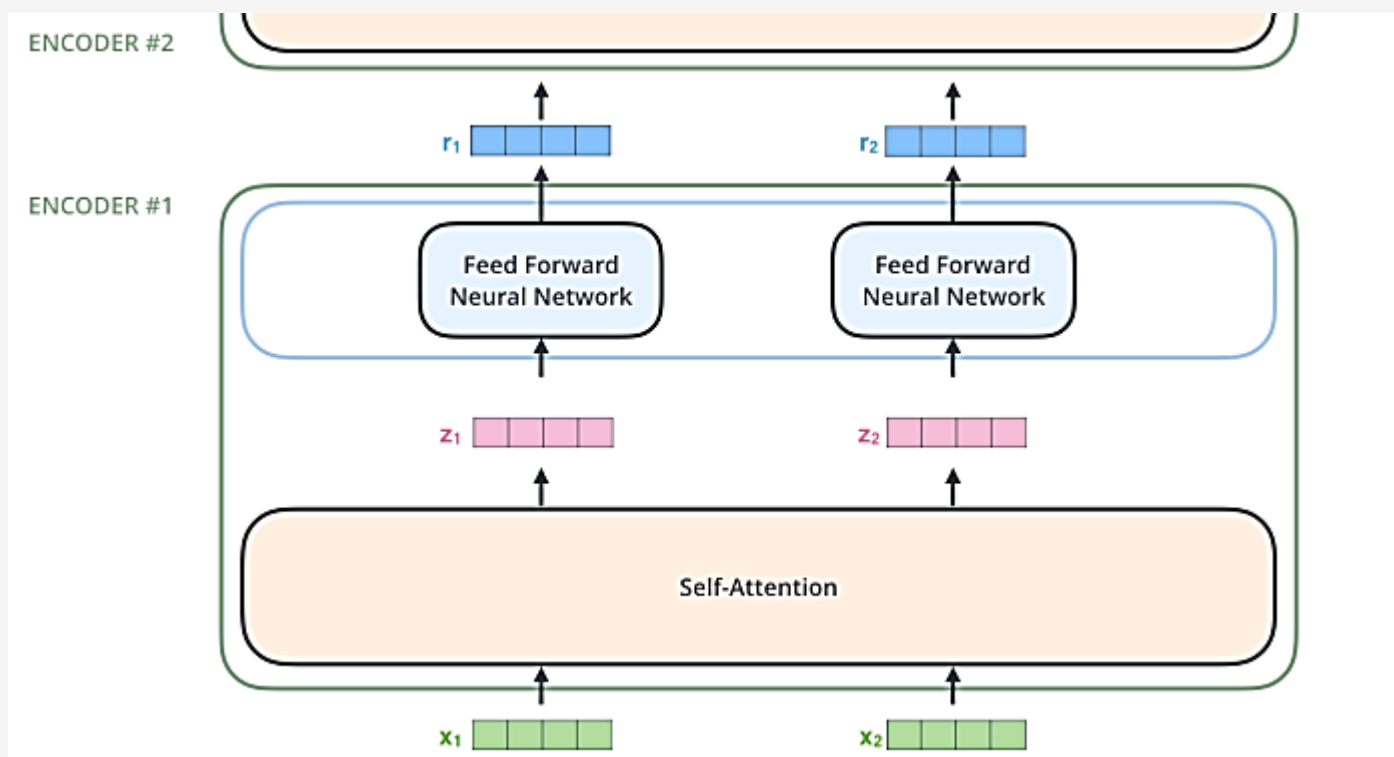
- دیکدر: خروجی انکدر که یک تعبیه است، دریافت می‌کند و مجدداً آن را به یک خروجی متنی تبدیل می‌کند؛
به بیان دیگر، نسخه ترجمه‌شده ورودی متنی را خروجی می‌دهد.

Transformers: بخش Encoder



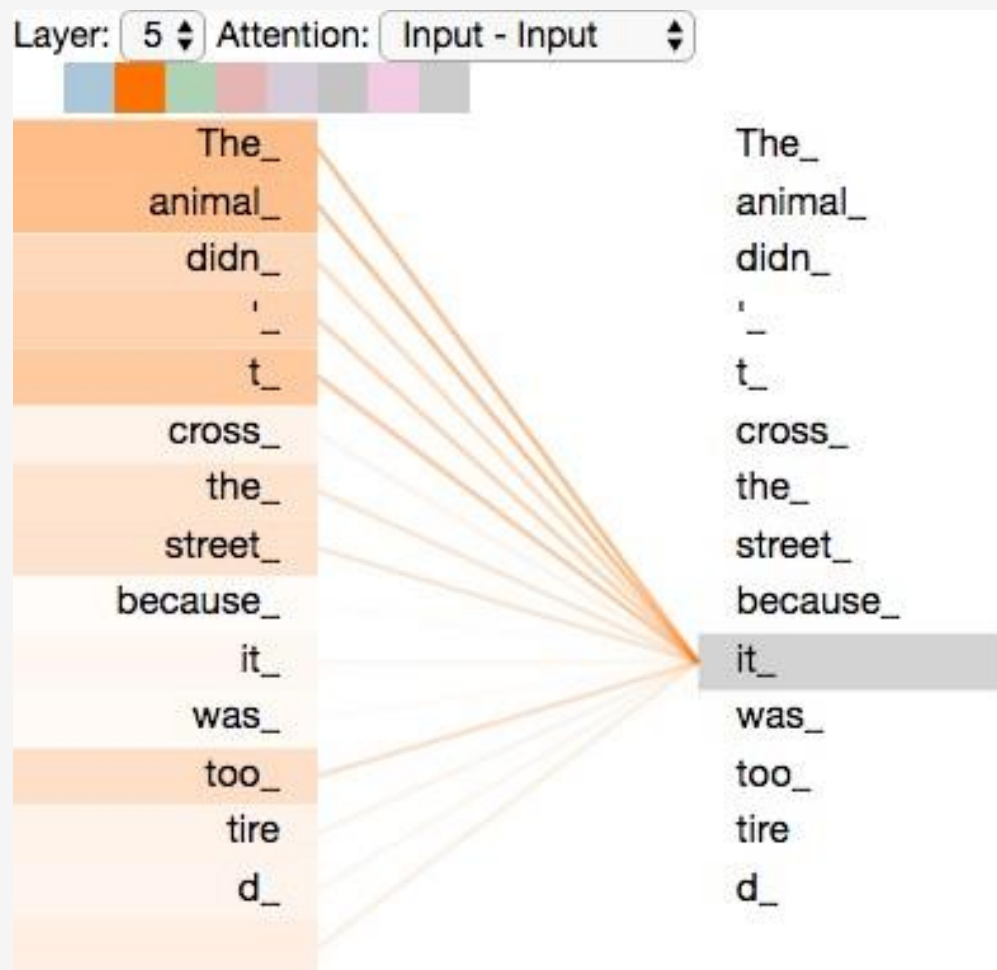
- تبدیل ورودی به بردار جاسازی
- انکودر بابت هر ورودی یک بردار z تولید می کند: ورودی لایه FF

Transformers



- نمایش لایه self-attention به صورت یکپارچه:
 - در نظر گرفتن وابستگی seq ورودی
- لایه FF جدا به ازای هر کلمه

Self Attention

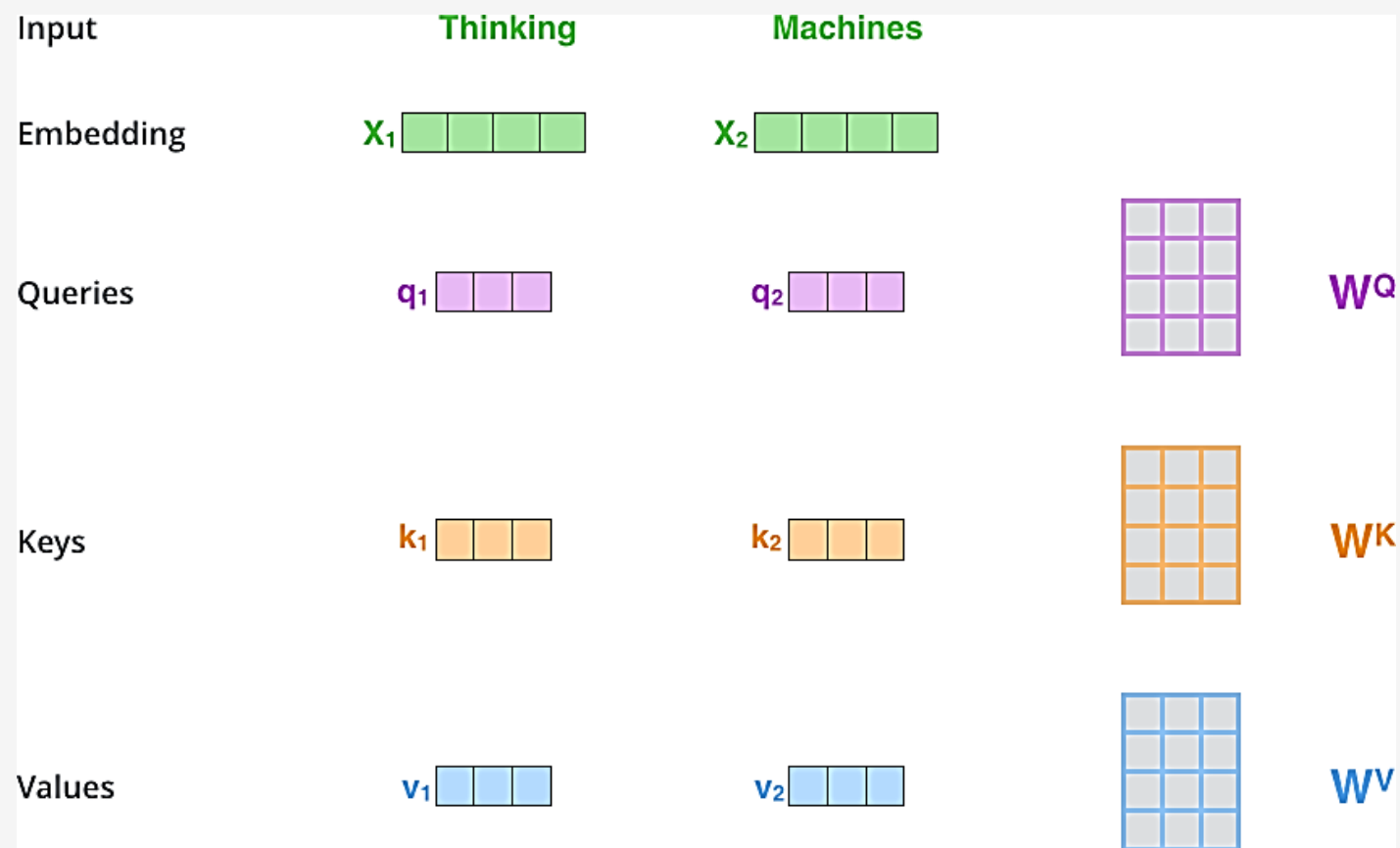


- به بخش‌های مختلف جمله نگاه می‌کند و تلاش می‌کند اطلاعات معنایی و نحوی بیشتری یاد بگیرد.
- پیاده سازی تاریخچه کلمات قبل و بعد که در بازگشتی ها بود.
- مشخص کردن وزن برای هر کلمه

The **animal** didn't cross the street because it was too **tired**.
The animal didn't cross the **street** because it was too **wide**.

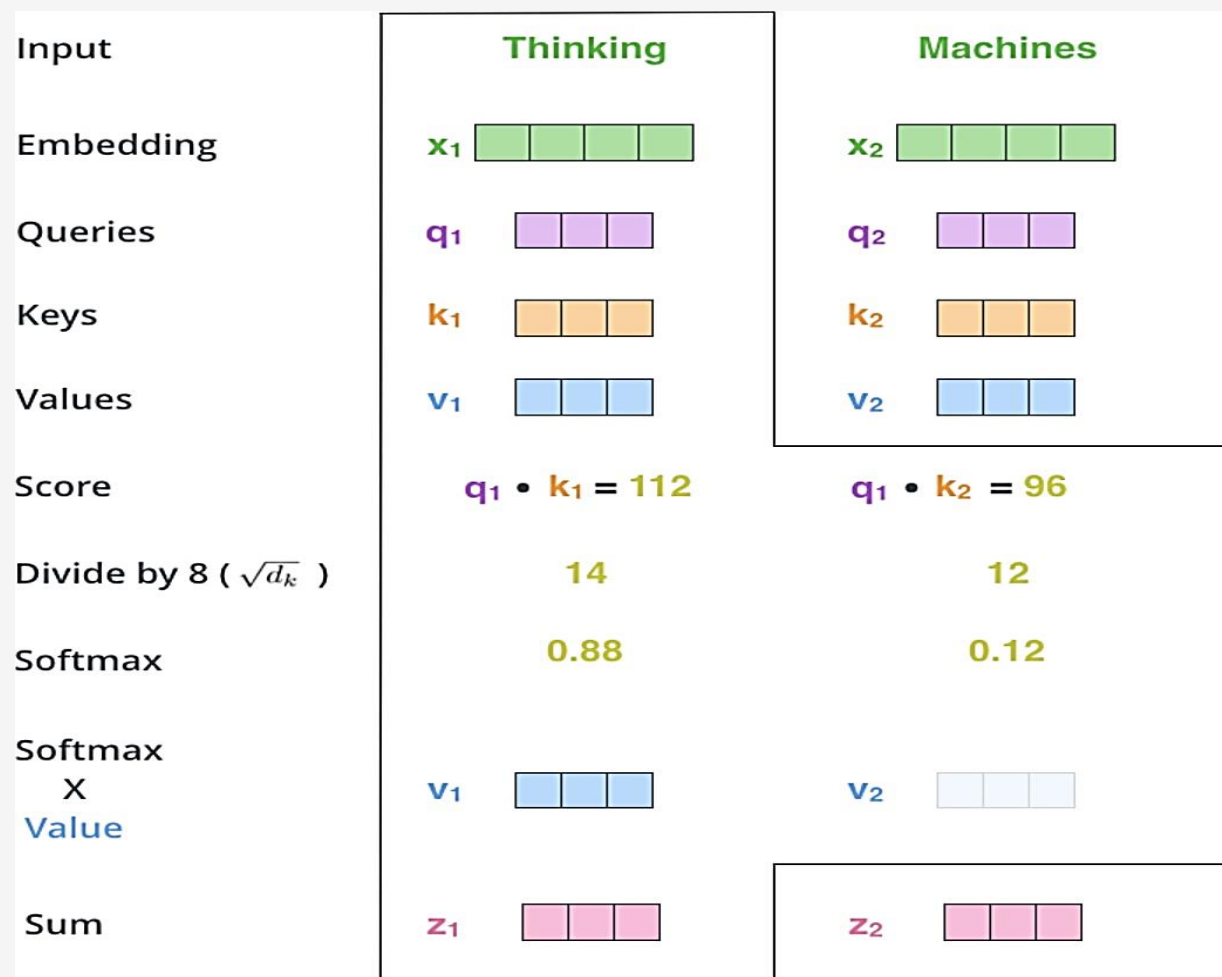
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Self Attention



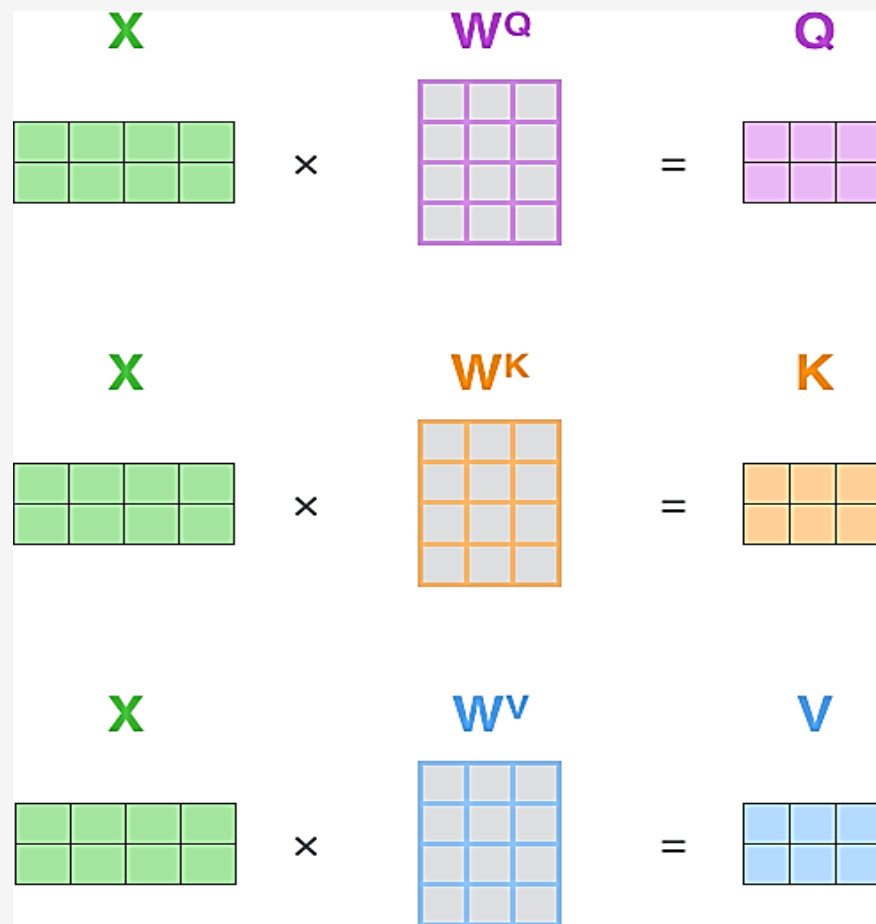
- در هر مرحله ۳ بردار Q ، K و V تشکیل می دهیم. چطور؟
- ۳ ماتریس وزن داریم، ضرب در بردار کلمه

Self Attention



- بعد از ورود هر کلمه به شبکه باید رتبه بندی باید داشته باشیم.
- چگونه امتیاز بدهیم؟
 - ماتریس q کلمه ضرب در k همه کلمات
 - با تقسیم بر رادیکال بعد k نرمالسازی
 - عبور از softmax (بازه ۰ و ۱ برای شناسایی مهم ها)
 - ضرب در بردار v کلمات دنباله
 - z به دست آمد.

Self Attention



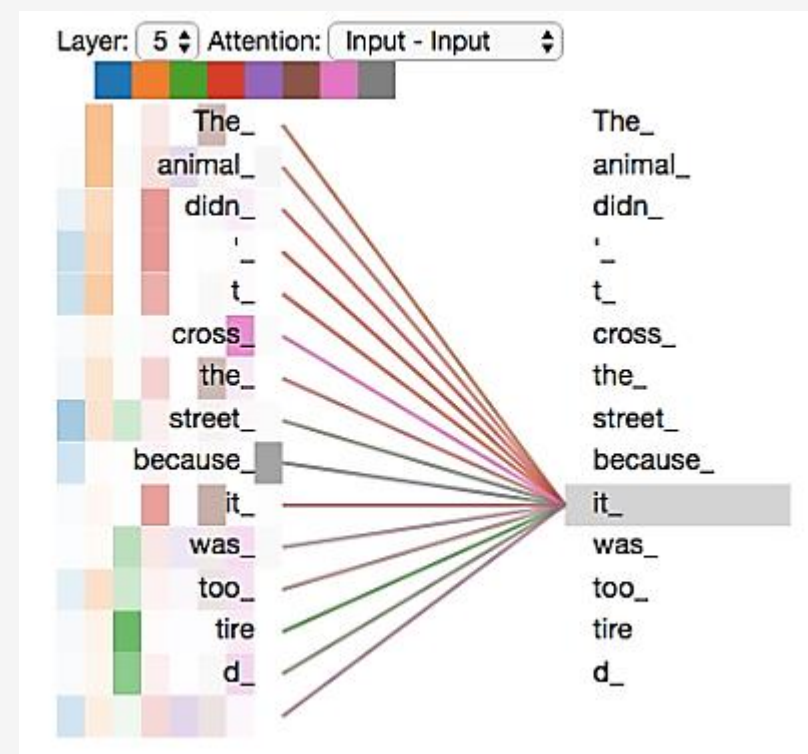
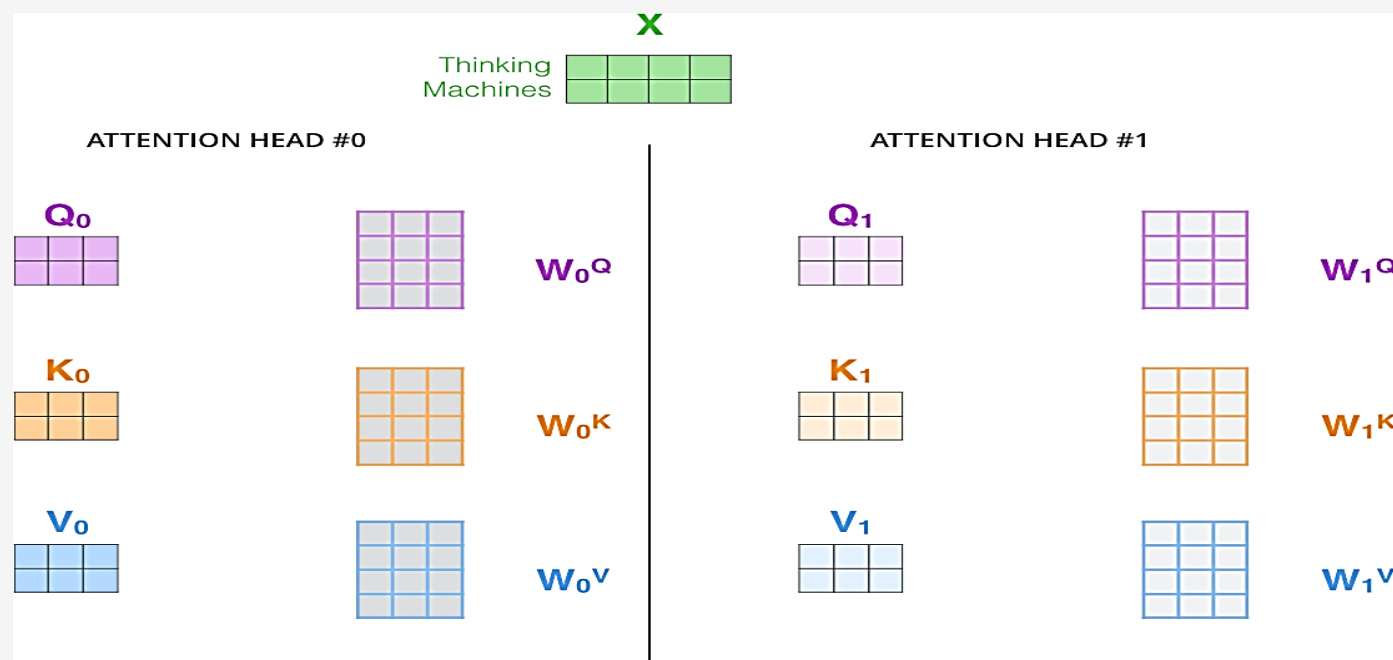
$$\text{softmax} \left(\frac{Q \times K^T}{\sqrt{d_k}} \right) \times V = Z$$

Diagram illustrating the Self Attention calculation:

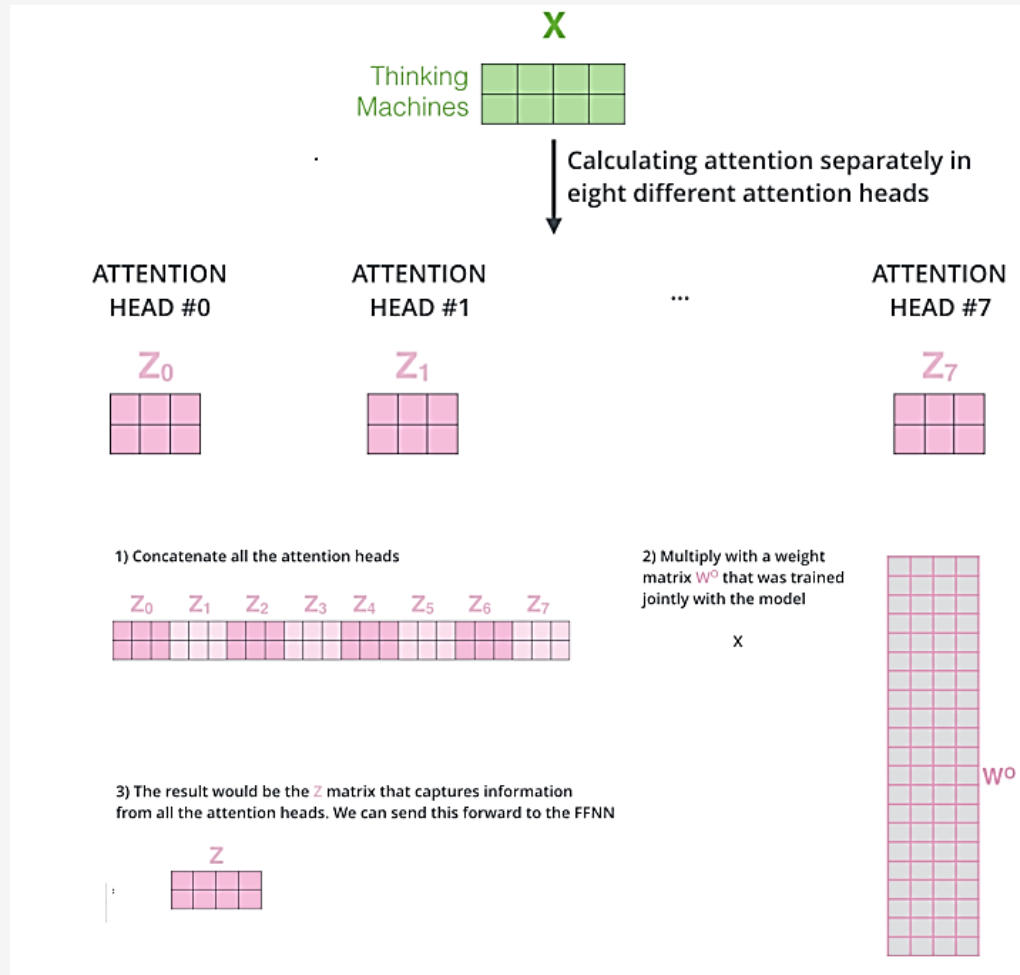
- Matrix Q (Purple) is multiplied by the transpose of matrix K (K^T , Orange).
- The result is divided by $\sqrt{d_k}$.
- The result is passed through a softmax function.
- The result is multiplied by matrix V (Blue) to produce the final output Z (Pink).

Multi Head Self Attention

- همزمان به چند جنبه اطلاعاتی و چند فضای بازنمایی توجه دارد.
- امکان موازی انجام دادن



Multi Head Self Attention



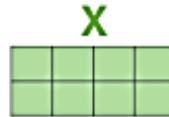
- تکرار همان کارها به ازای هر هد و بدست آوردن چند Z
- تبدیل چند Z به یک Z مثل قبل
- وصل کردن Z ها به هم و ضرب در ماتریس مناسب

Multi Head Self Attention

1) This is our input sentence*

Thinking
Machines

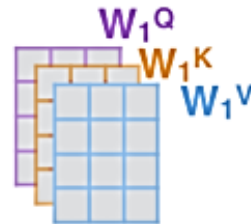
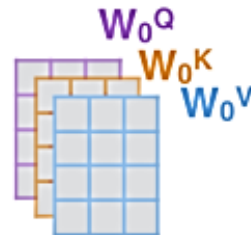
2) We embed each word*



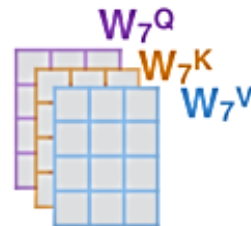
* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one



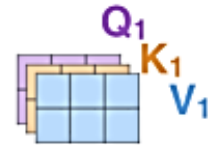
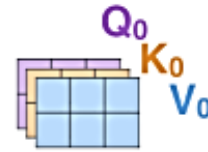
3) Split into 8 heads. We multiply X or R with weight matrices



...



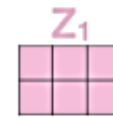
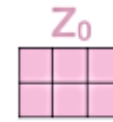
4) Calculate attention using the resulting $Q/K/V$ matrices



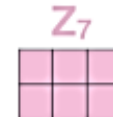
...



5) Concatenate the resulting Z matrices, then multiply with weight matrix W^O to produce the output of the layer



...



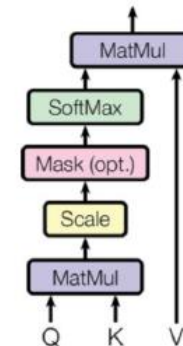
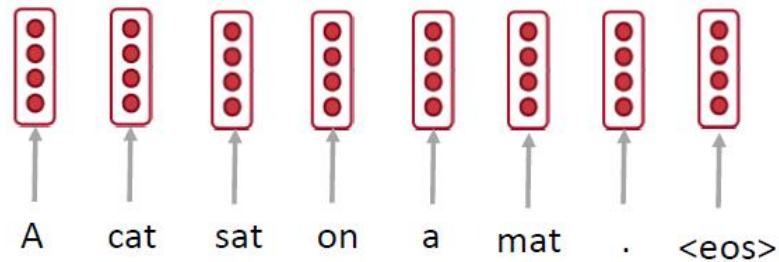
خروجی لایه قبل



Self-attention: A Running Example

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

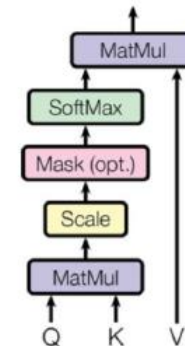
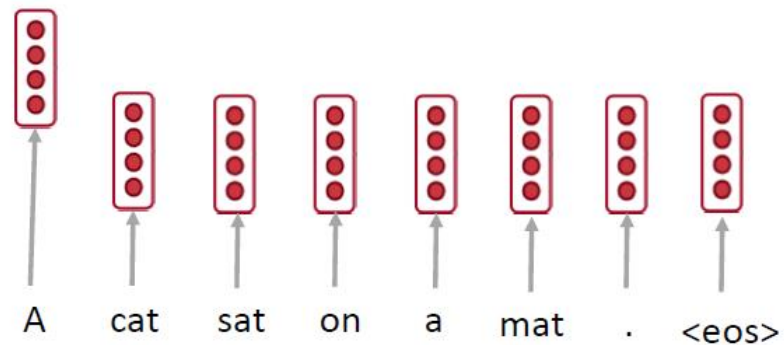
Scaled Dot-Product Attention



Self-attention: A Running Example

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

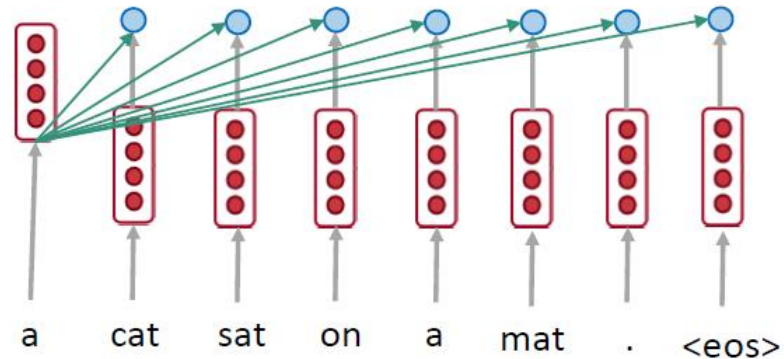
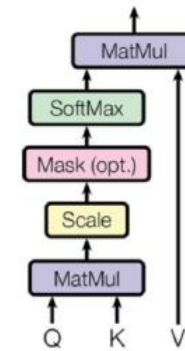
Scaled Dot-Product Attention



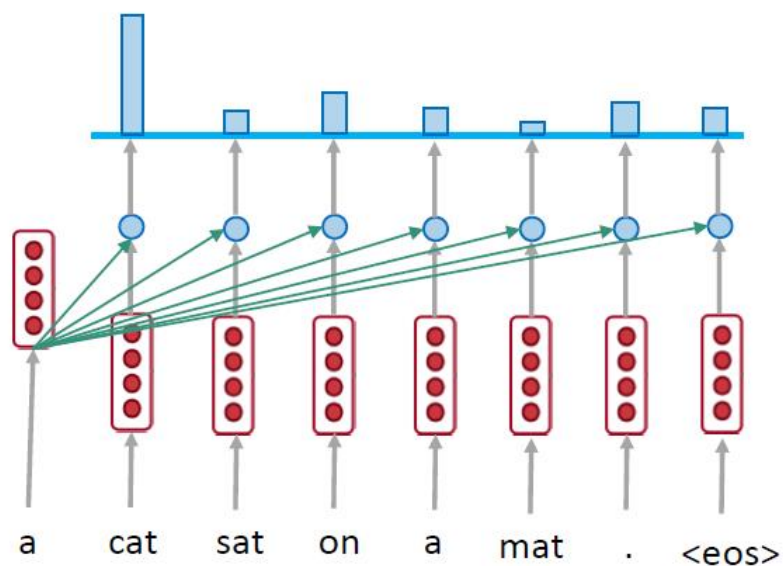
Self-attention: A Running Example

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

Scaled Dot-Product Attention

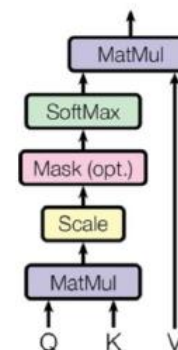


Self-attention: A Running Example

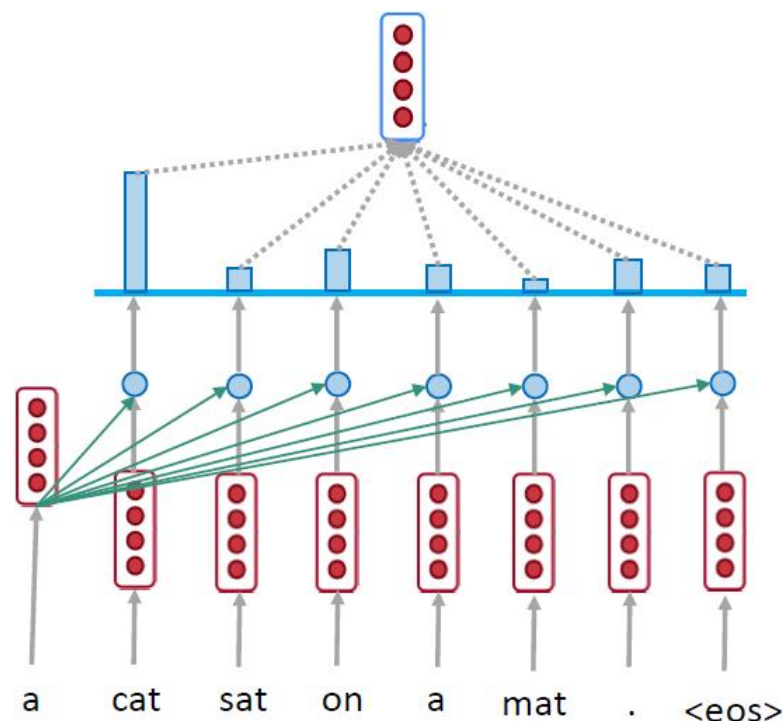


$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Scaled Dot-Product Attention

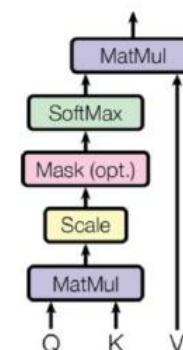


Self-attention: A Running Example



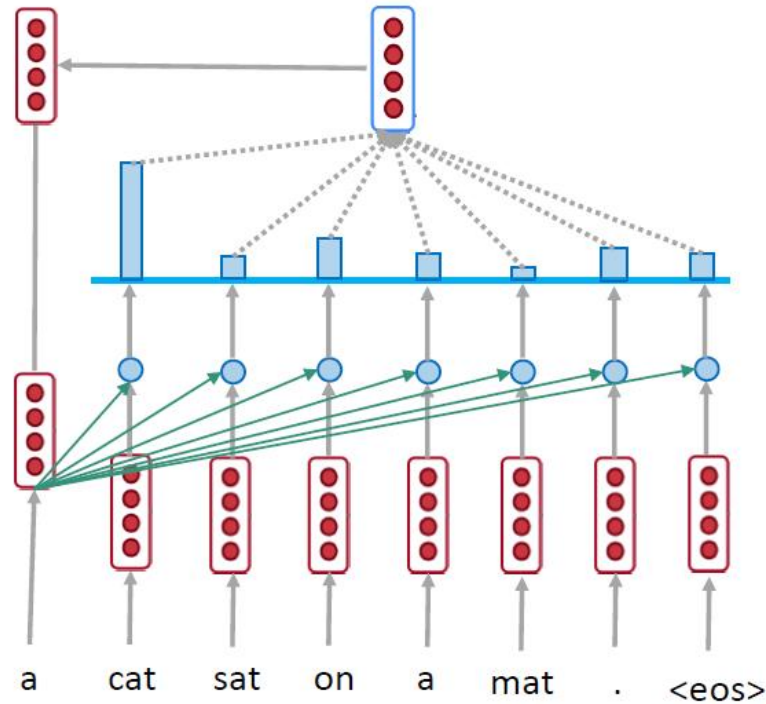
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Scaled Dot-Product Attention



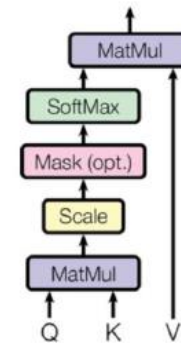
Self-attention: A Running Example

update
representation
for the word "a"

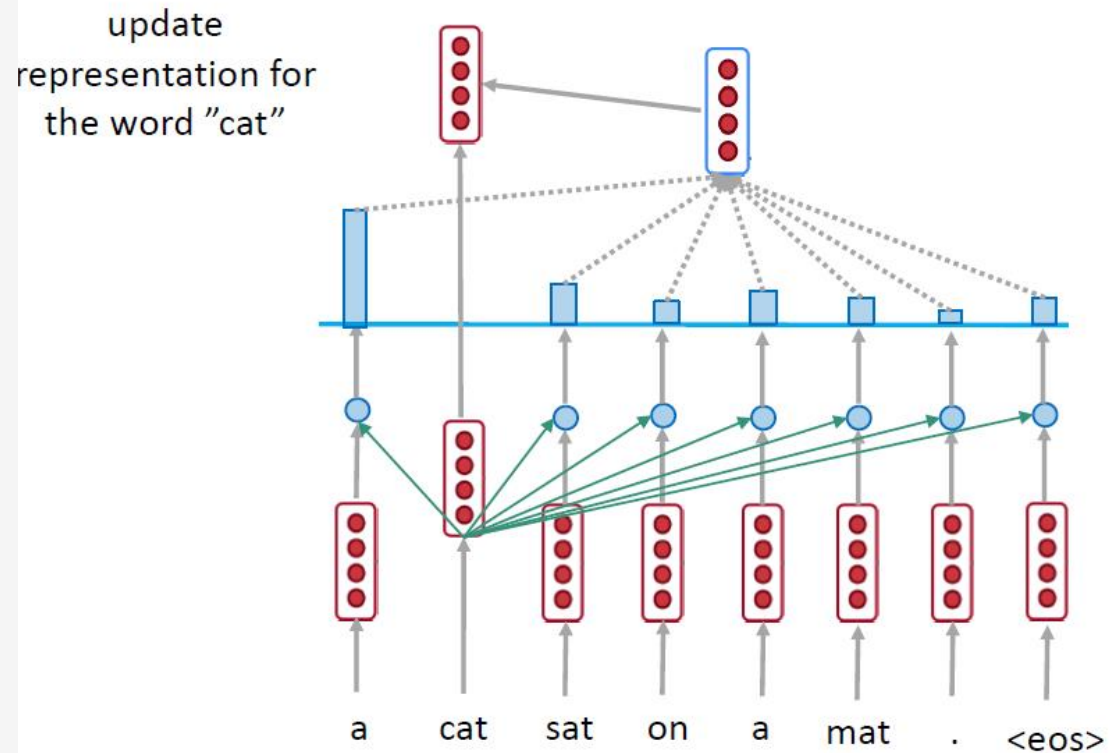


$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Scaled Dot-Product Attention

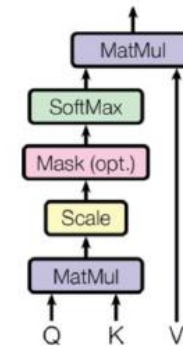


Self-attention: A Running Example



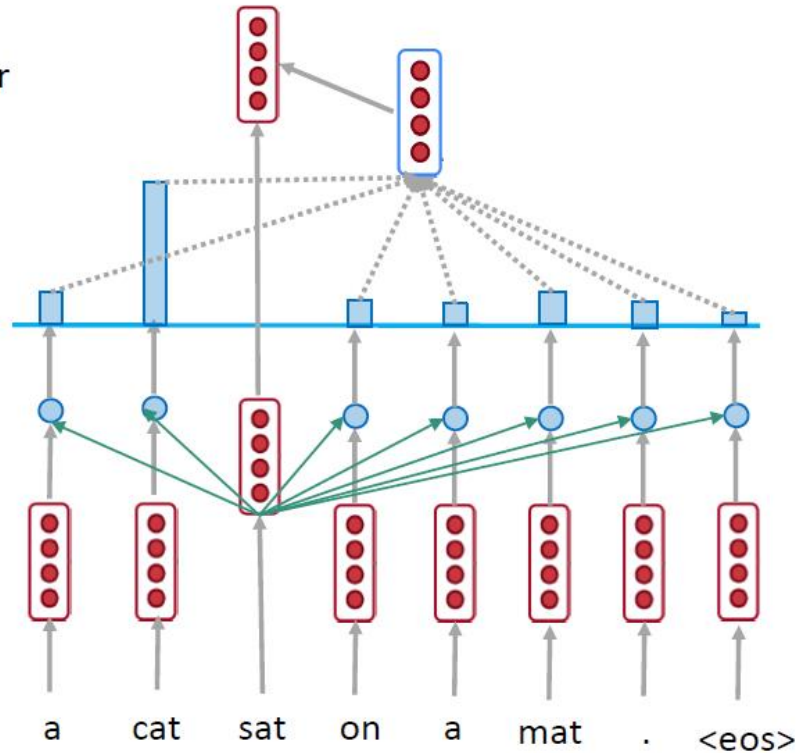
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Scaled Dot-Product Attention



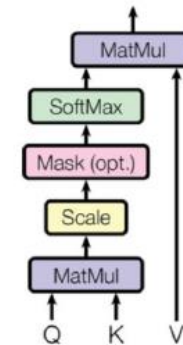
Self-attention: A Running Example

update
representation for
the word "sat"



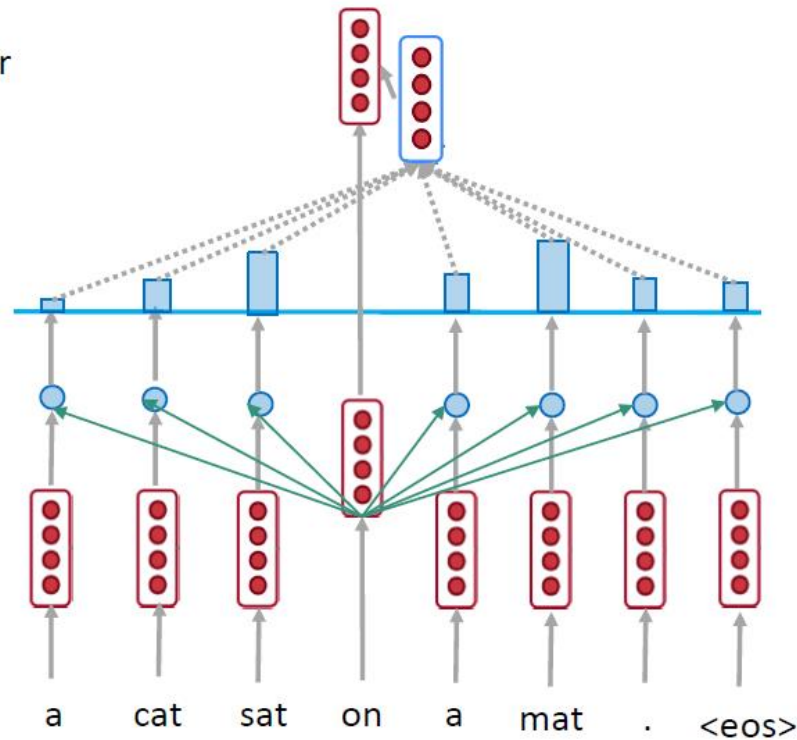
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Scaled Dot-Product Attention



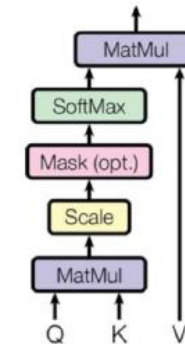
Self-attention: A Running Example

update
representation for
the word "on"



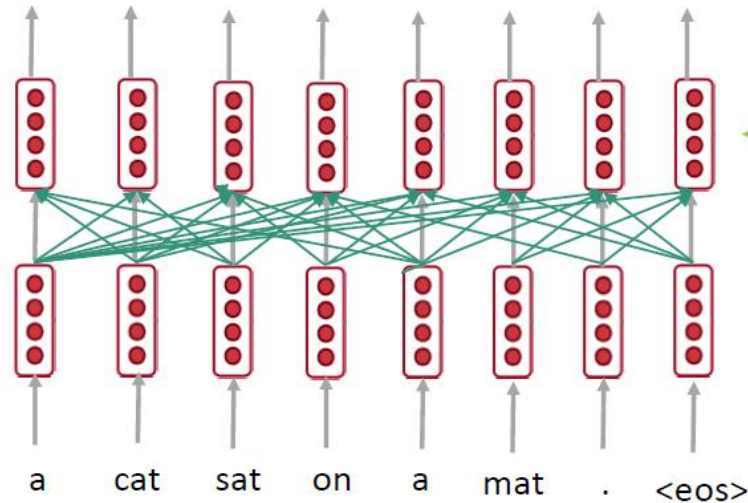
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Scaled Dot-Product Attention



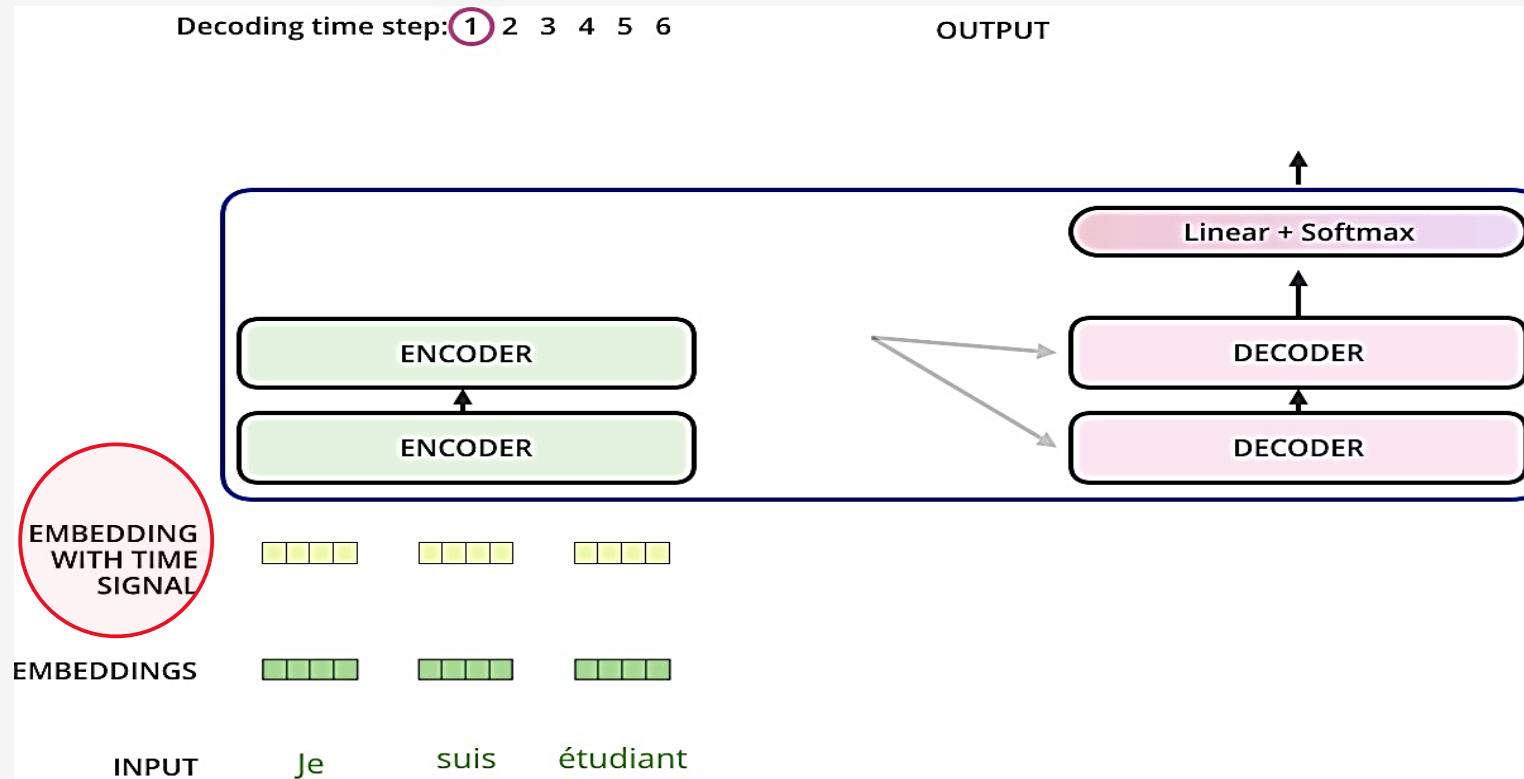
Self-attention: A Running Example

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

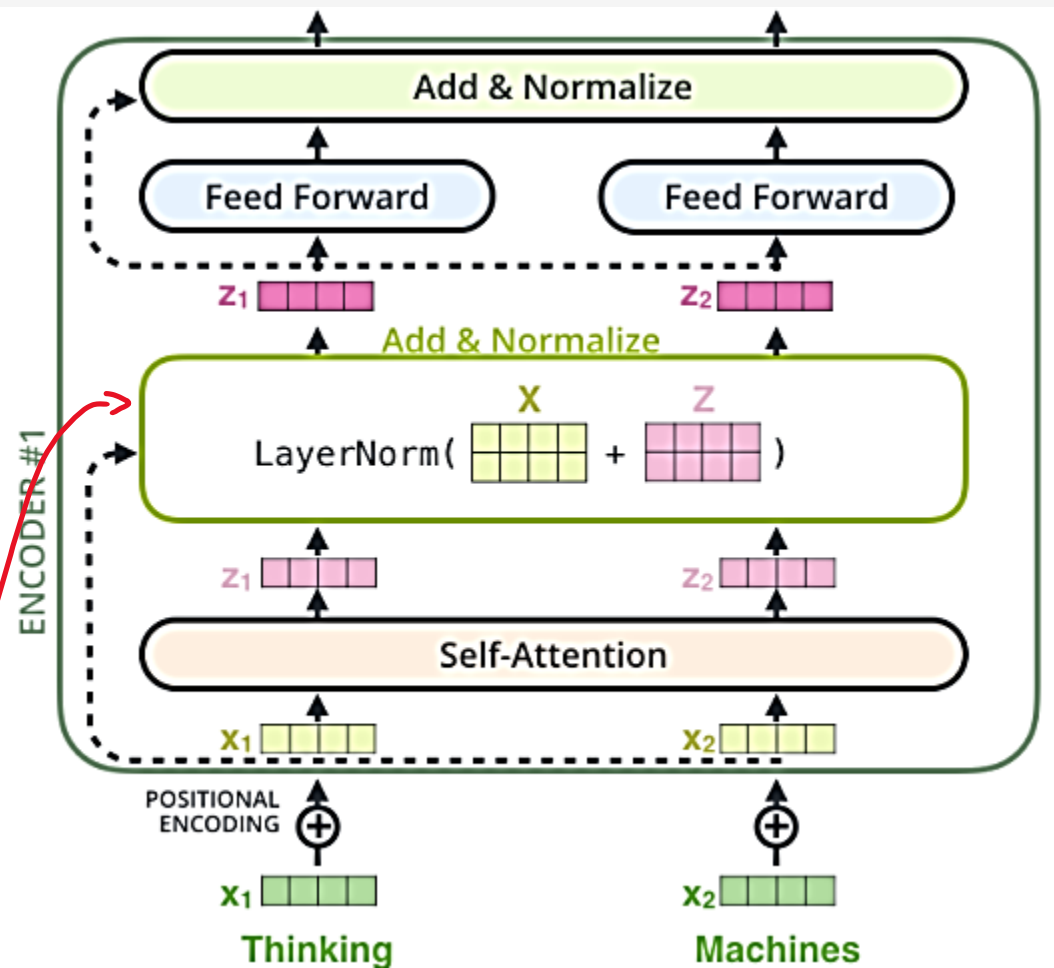
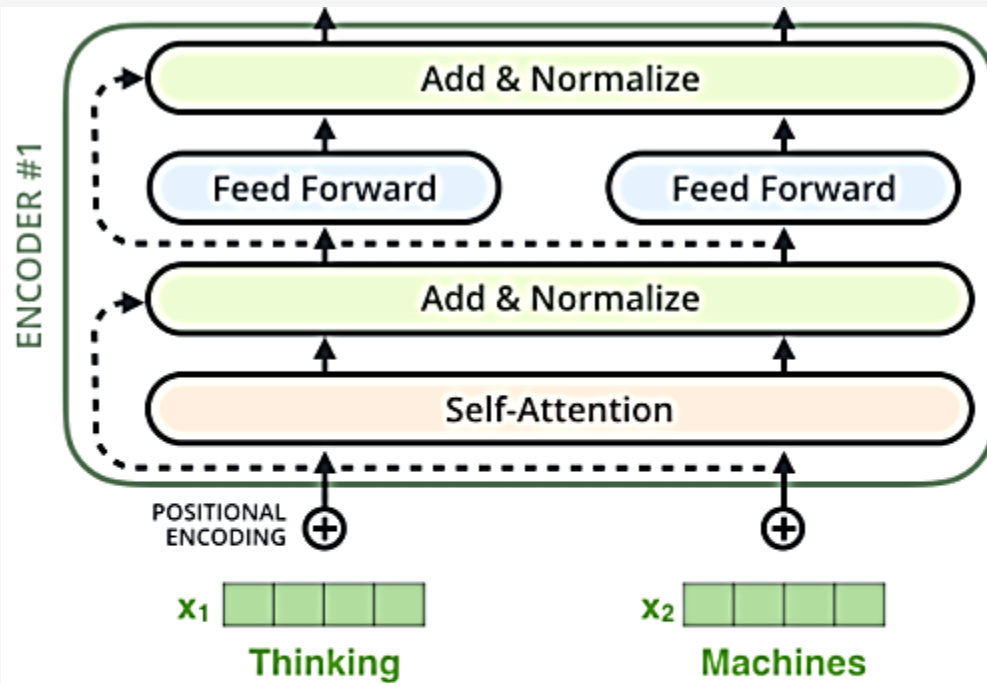


Updated representations
for each word
(one layer)

پیاده سازی توالی

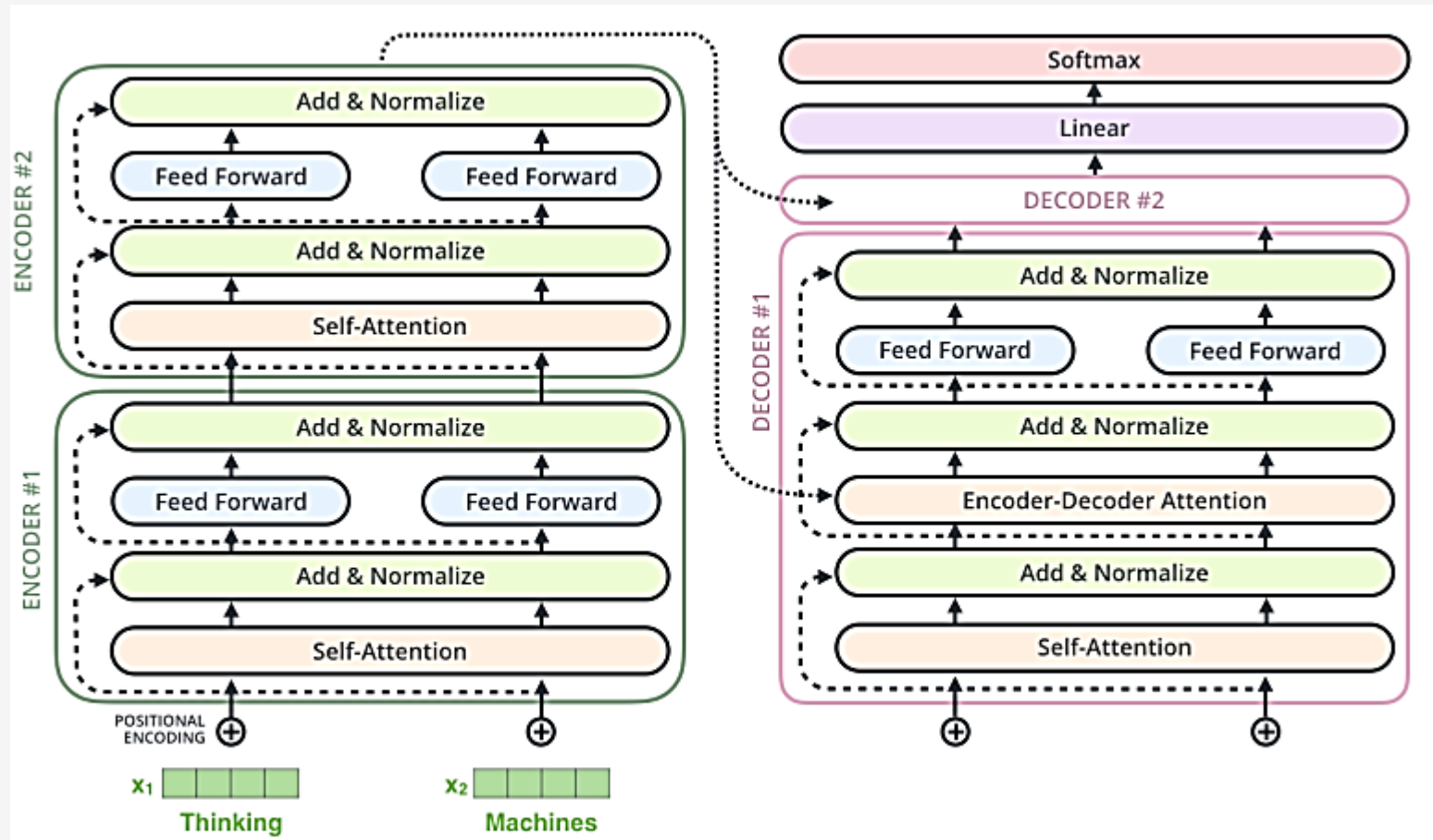


Transformers: بخش Add & Normalize



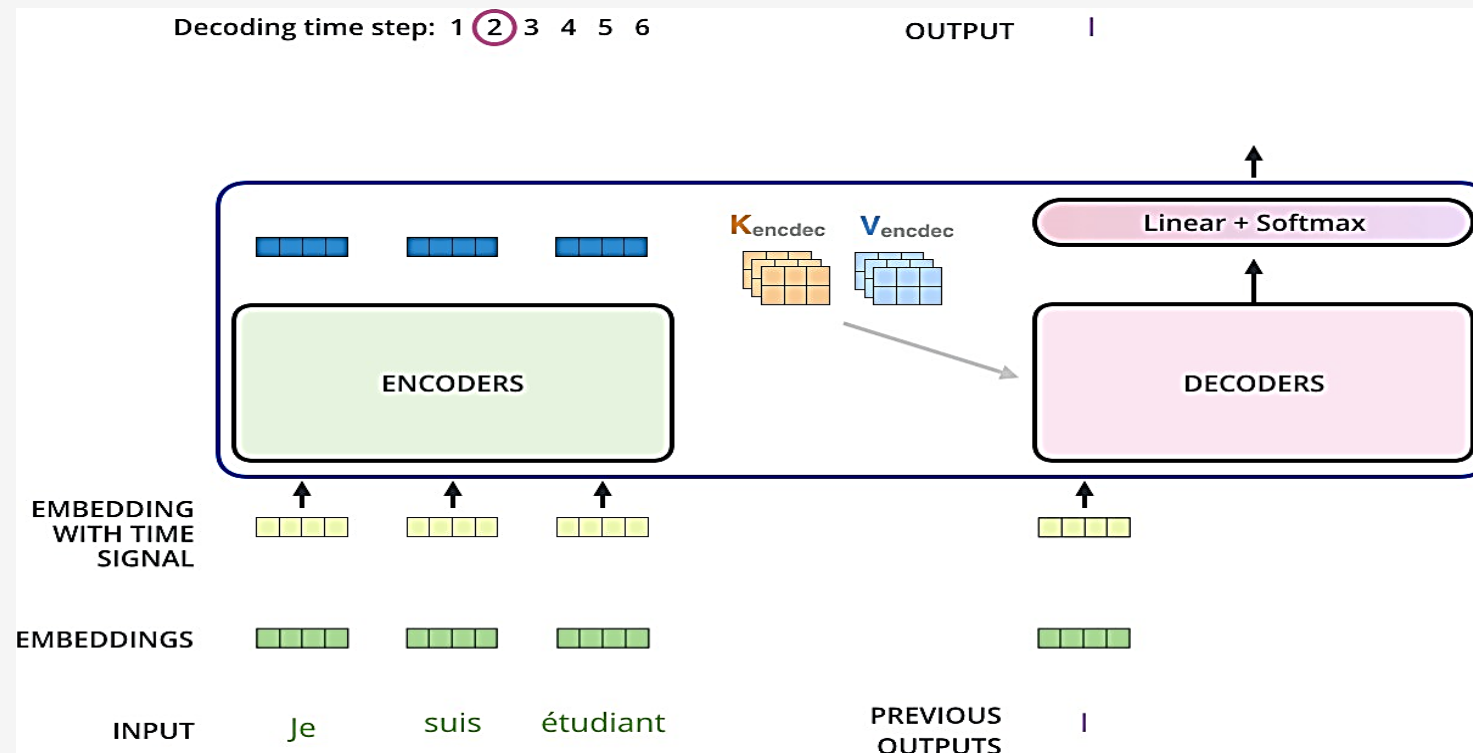
- برای افزایش سرعت آموزش
- عبور ماتریس جمع X و Z از یک تابع نرمال

Transformers: بخش Add & Normalize



Transformers: بخش Decoder

- تا اینجا دیدیم ورودی بعد از لایه توجه نرمال شده و به FF می رود.
- خروجی FF ورودی دیکودر است.
- دیکودر مشابه انکودر است، فقط یک لایه encoder-decoder attention بیشتر دارد.
- همان Multi head attention، فقط ماتریس های key و Value را از انکودر می گیرد.



Transformers: لایه softmax

- مشابه قبلی ها

Which word in our vocabulary
is associated with this index?

Get the index of the cell
with the highest value
(argmax)

log_probs



am

5

Softmax

logits

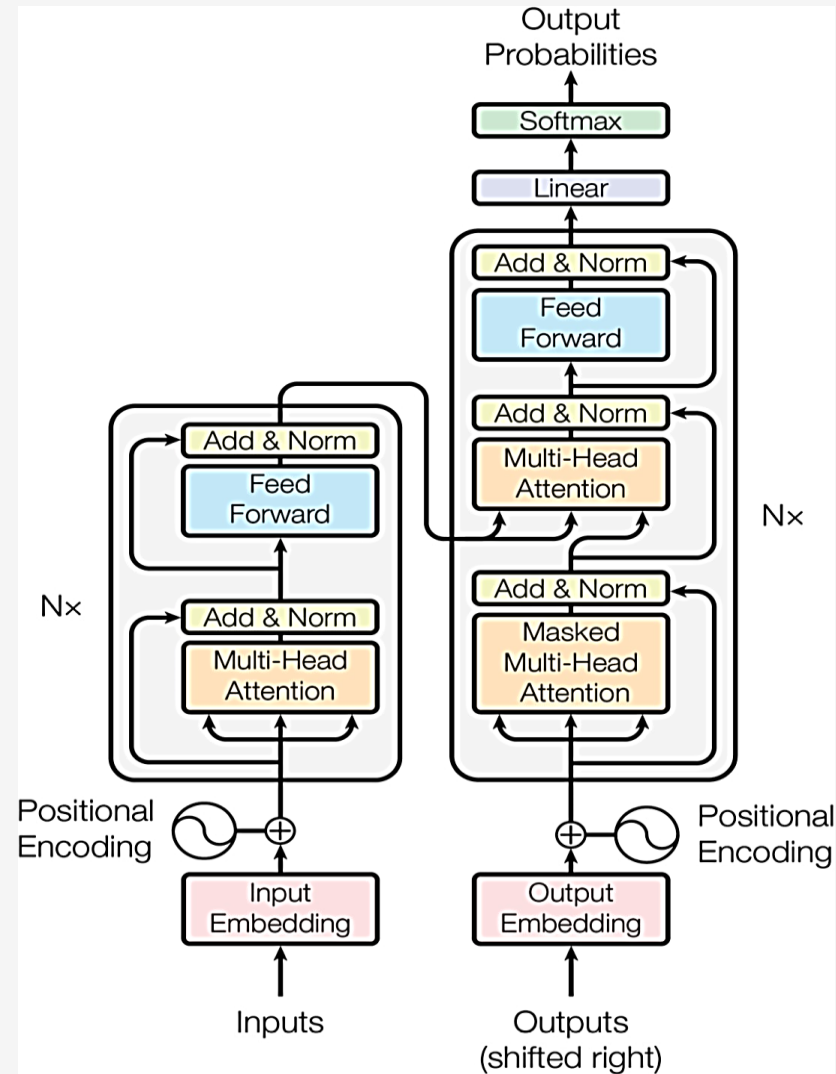


Linear

Decoder stack output



Transformers: معماری کلی

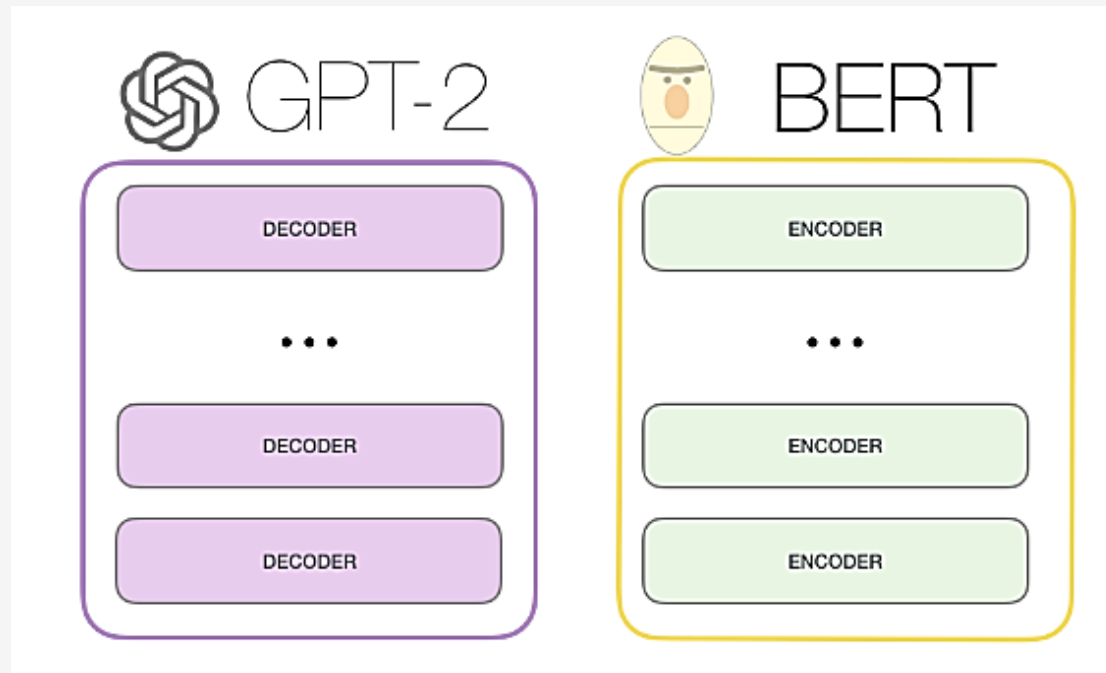


(Vaswani *et al.*, 2017)

Transformer Networks Visualization



GPT و BERT



- یک ترنسفورمر از پشته انکودر برای مدل سازی ورودی استفاده می کند و از پشته دیکودر برای مدل سازی خروجی (با استفاده از اطلاعات ورودی از سمت رمزگذار) اگر ورودی نداریم، فقط می خواهیم «کلمه بعدی» را مدل کنیم، می توانیم سمت رمزگذار یک ترنسفورمر را حذف کنیم و «کلمه بعدی» را یکی یکی خروجی بدهیم.
- این به ما GPT می دهد.
- اگر ما فقط علاقه مند به آموزش یک مدل زبانی برای ورودی برای برخی از تسک های دیگر هستیم، به رمزگشای ترنسفورمر نیاز نداریم که به ما BERT می دهد.