

Attention is all you need

RNN is unnecessary

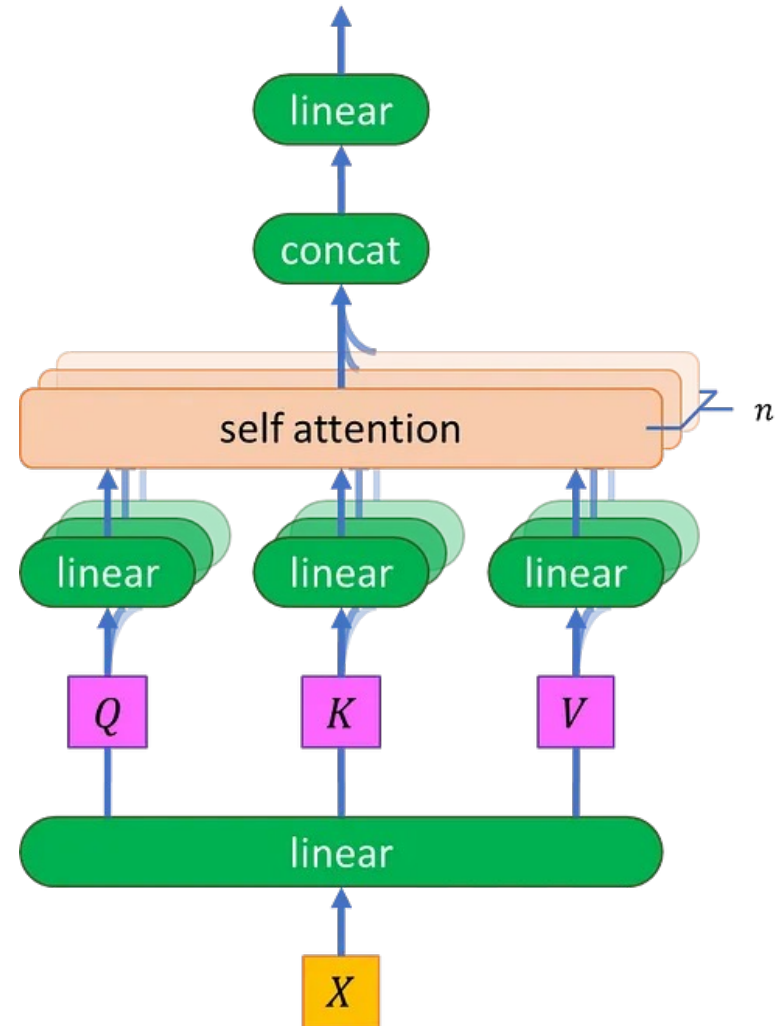
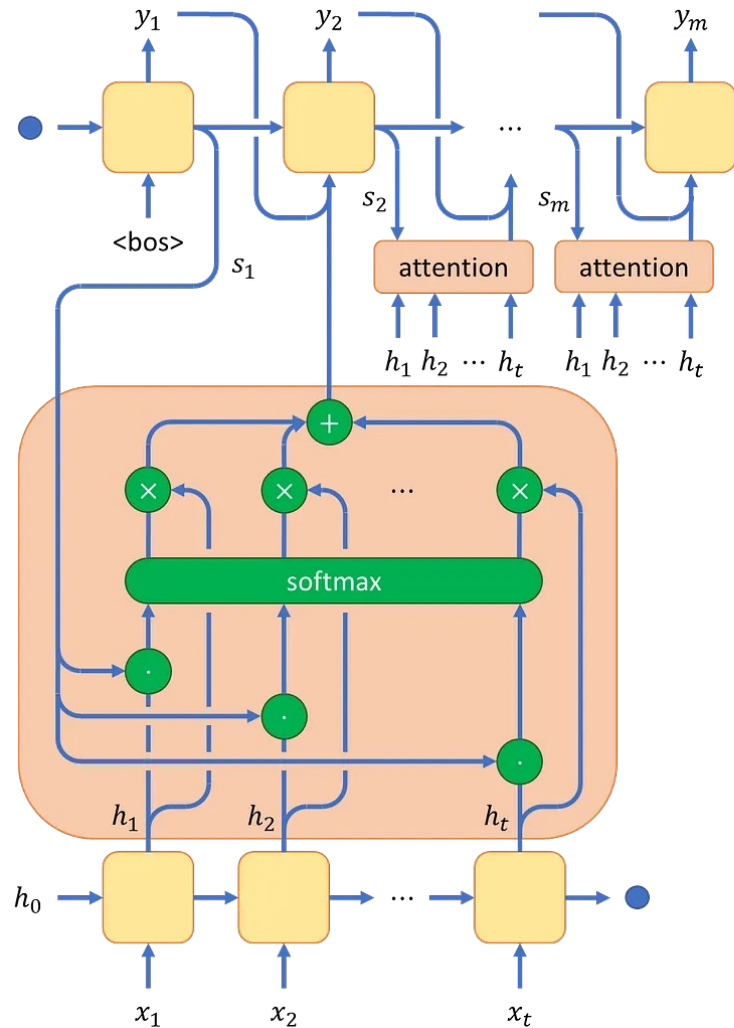
Liang Hong

17 Mar, 2023

Outline

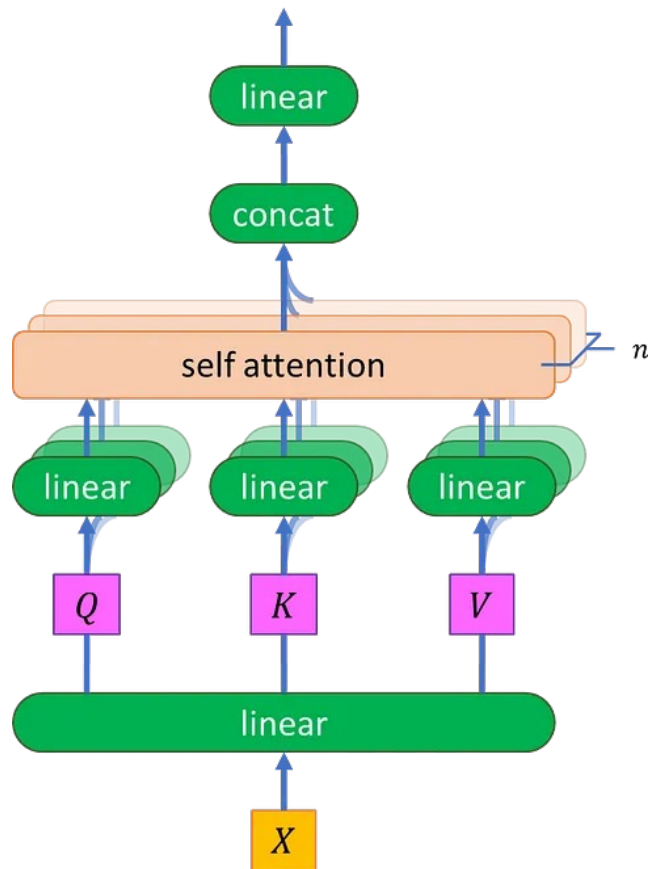
- How inference change from RNN style to Transformer style
- How attention deal with token order
- What's next

Inference style change



Token order

- Without sequential input, how token order is perceived?



Attention is performed with dot product, and the previous linear layers are shared.

-> permutation irrelevant

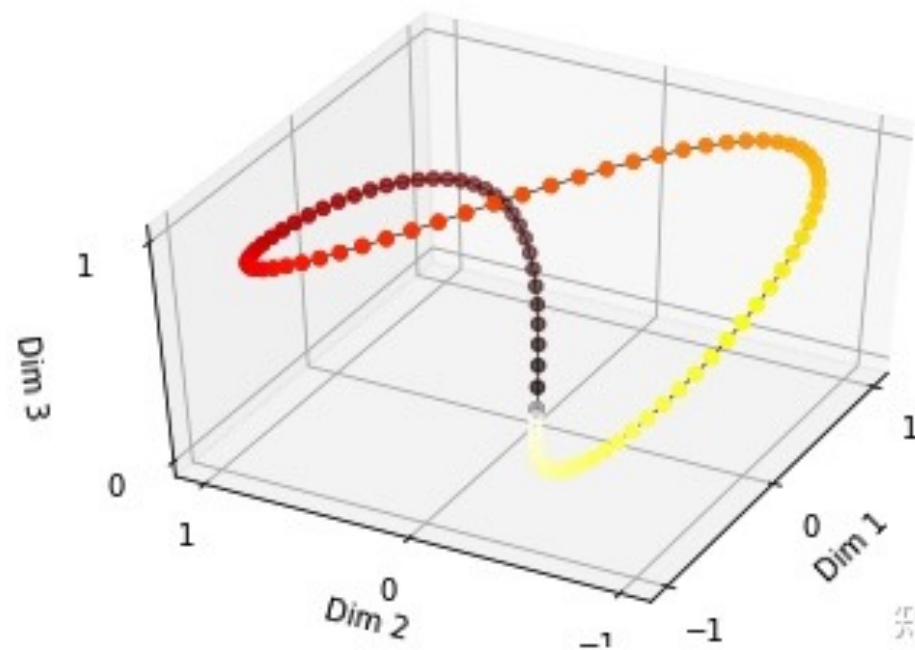
Positional encoding

$$\begin{cases} PE(pos, 2i) = \sin(pos/10000^{2i/d_{model}}) \\ PE(pos, 2i+1) = \cos(pos/10000^{2i/d_{model}}) \end{cases}$$

Why this frequency

- ✗ int type
- ✗ 0~1 representation

Credit to [@lemonround](#)



知乎 @猛犸

Positional encoding

$$\begin{cases} PE(pos, 2i) = \sin\left(pos/10000^{2i/d_{model}}\right) \\ PE(pos, 2i + 1) = \cos\left(pos/10000^{2i/d_{model}}\right) \end{cases}$$

(now we have unique and continuous position rep vectors)

How can we get relative position from position encoding?

$$\begin{aligned} PE_t^T * PE_{t+\Delta t} &= \sum_{i=0}^{\frac{d_{model}}{2}-1} [\sin(w_i t) \sin(w_i (t + \Delta t)) + \cos(w_i t) \cos(w_i (t + \Delta t))] \\ &= \sum_{i=0}^{\frac{d_{model}}{2}-1} \cos(w_i (t - (t + \Delta t))) \\ &= \sum_{i=0}^{\frac{d_{model}}{2}-1} \cos(w_i \Delta t) \end{aligned}$$

$$PE_t^{(i)} = \begin{cases} \sin(w_i t), & \text{if } k = 2i \\ \cos(w_i t), & \text{if } k = 2i + 1 \end{cases}$$

Positional encoding

$$\begin{cases} PE(pos, 2i) = \sin(pos/10000^{2i/d_{model}}) \\ PE(pos, 2i+1) = \cos(pos/10000^{2i/d_{model}}) \end{cases}$$

(now we can get relative positions from positional encoding)

$$\begin{aligned} PE_t^T * PE_{t+\Delta t} &= \sum_{i=0}^{\frac{d_{model}}{2}-1} [\sin(w_i t) \sin(w_i (t + \Delta t)) + \cos(w_i t) \cos(w_i (t + \Delta t))] \\ &= \sum_{i=0}^{\frac{d_{model}}{2}-1} \cos(w_i (t - (t + \Delta t))) \\ &= \sum_{i=0}^{\frac{d_{model}}{2}-1} \cos(w_i \Delta t) \end{aligned}$$



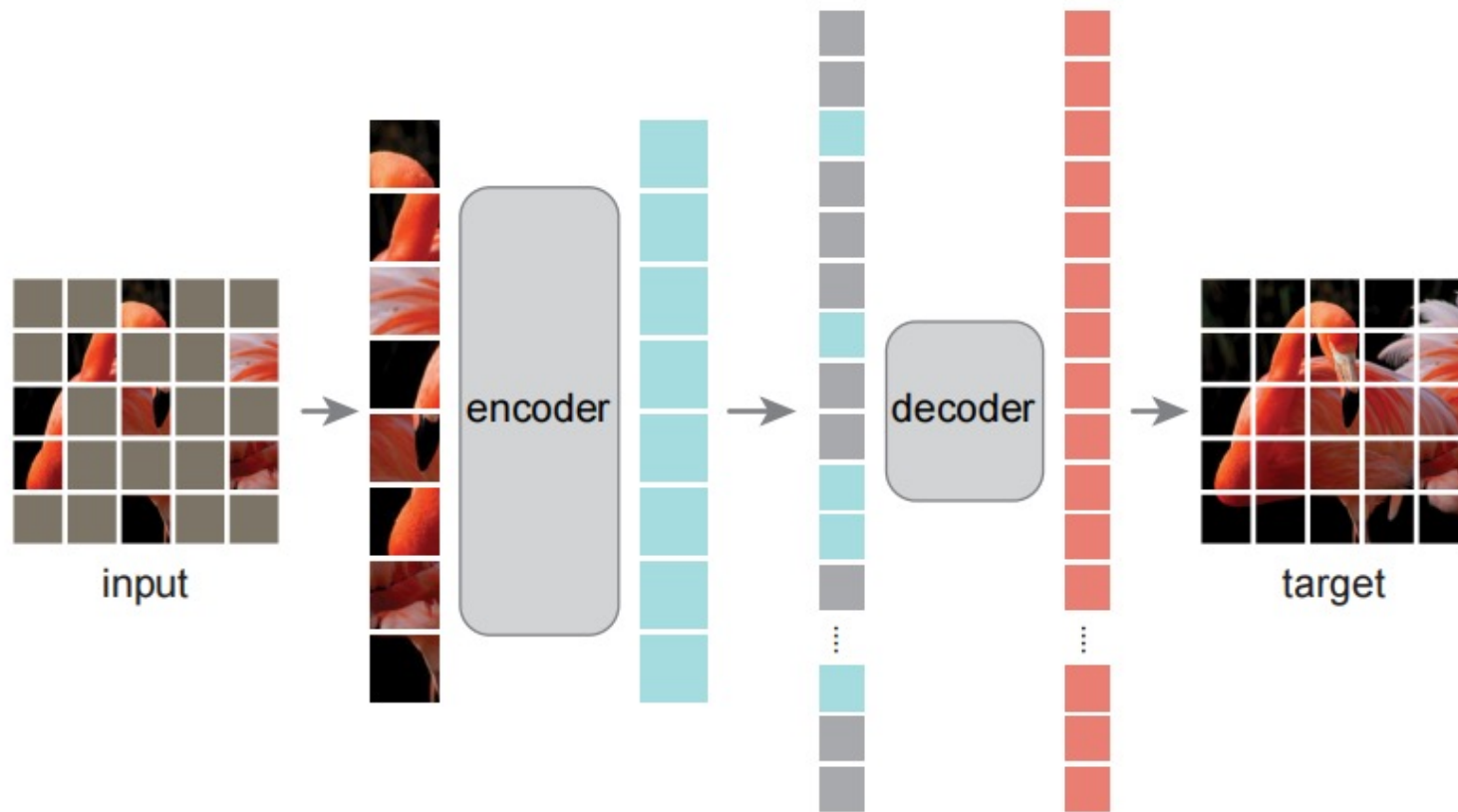
distance aware



direction aware

$$PE_t^{(i)} = \begin{cases} \sin(w_i t), & \text{if } k = 2i \\ \cos(w_i t), & \text{if } k = 2i + 1 \end{cases}$$

What's next



- MAE: Masked Autoencoders Are Scalable Vision Learners

What's next

