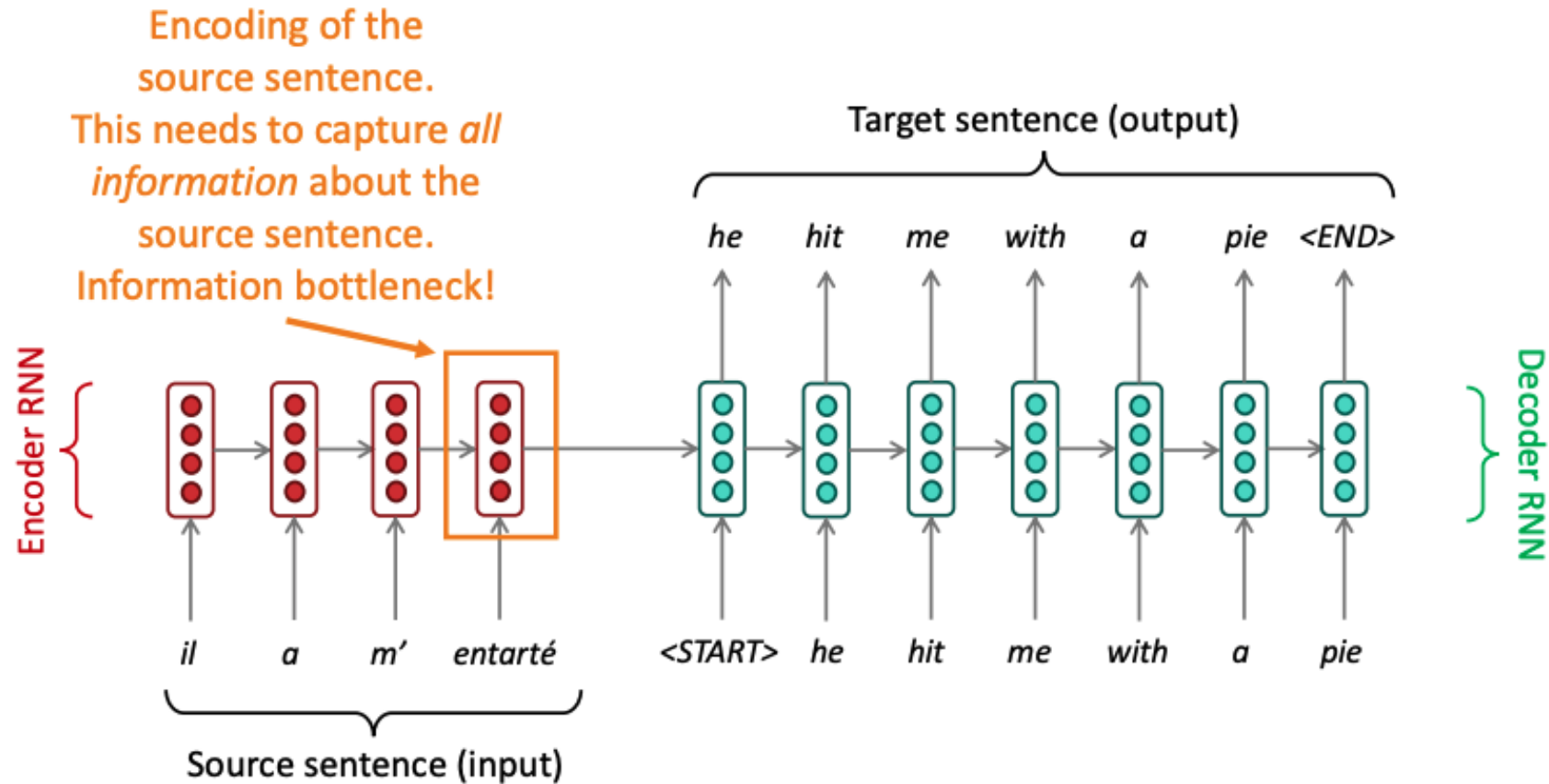


transformer

stats403_deep_learning
spring_2025
lecture_7

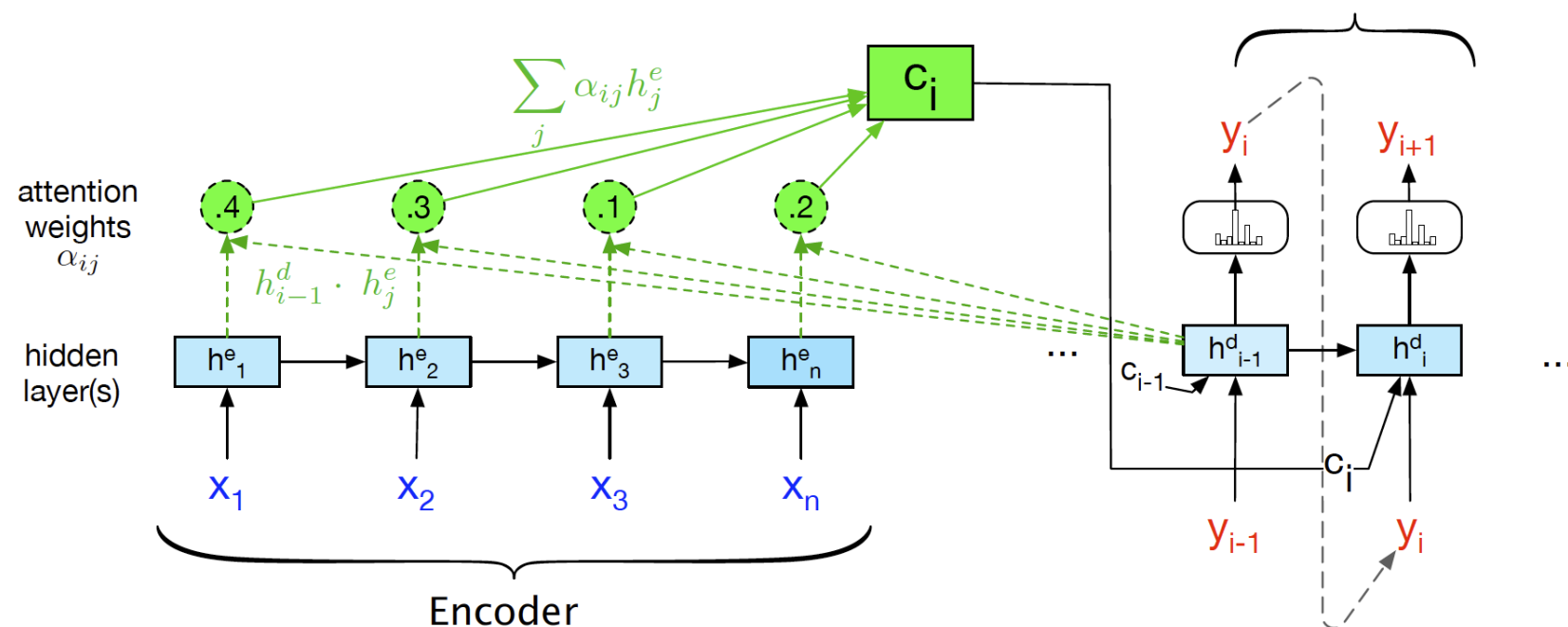
7.1 attention and self-attention

drawback of RNN encoder-decoder



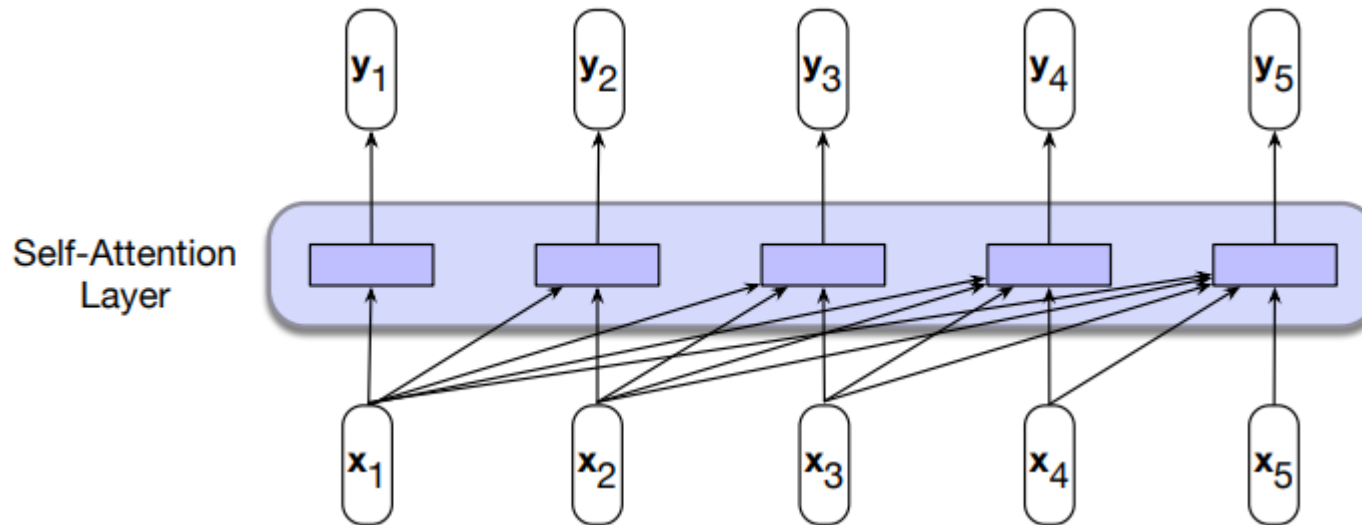
machine translation: bottleneck

- idea: map all hidden states to all output state by learning an alignment with weights
- use “attention” to determine the importance



transformer

- a transformer maps a sequence of input vectors $(\mathbf{x}_1, \dots, \mathbf{x}_n)$ to a sequence of output vectors $(\mathbf{y}_1, \dots, \mathbf{y}_n)$
- the most important layer in a transformer is the **self-attention layer**



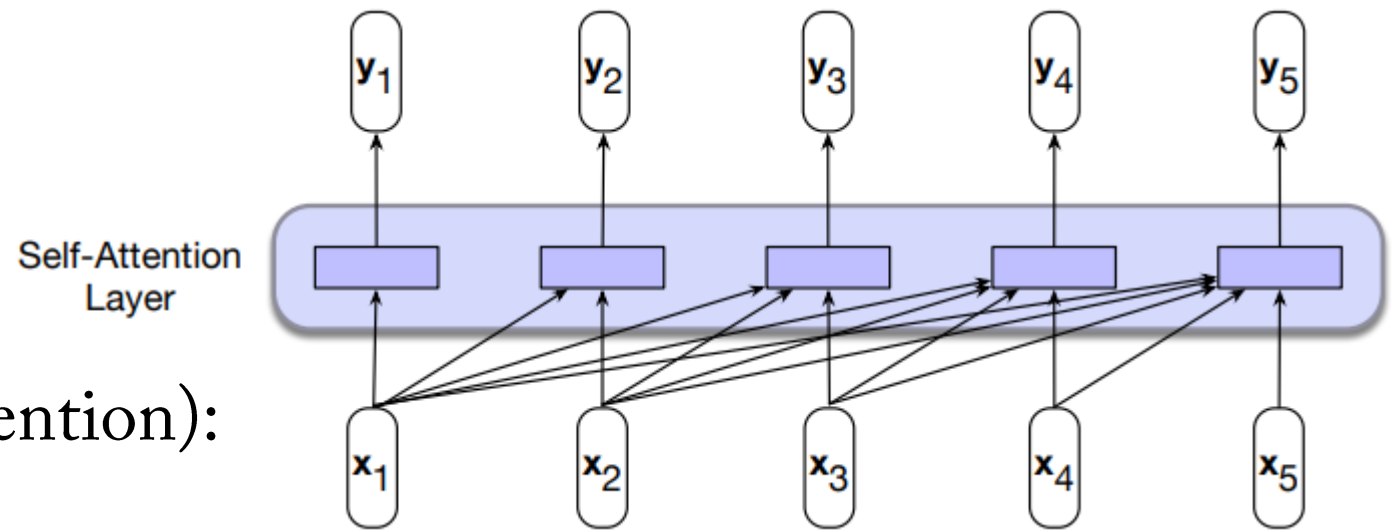
an example of (causal masked) self-attention which attends to all the inputs up to the current one

self-attention

- **attention**: compare an item of interest to a collection of other items in a way that reveals their relevance in the current context
- **self-attention**: the comparison is done within a given sequence
- the result of comparison is used to compute the output

self-attention

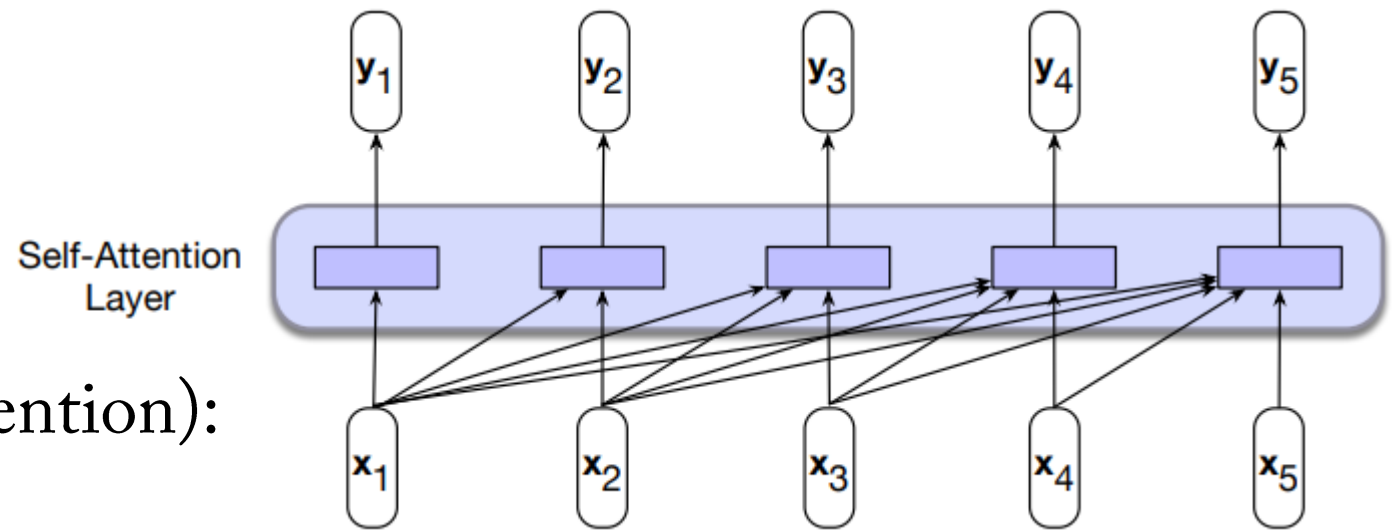
- example (a simple self-attention):



- the simplest form of comparison is a dot product, e.g.,
$$\text{score}(\mathbf{x}_i, \mathbf{x}_j) = \langle \mathbf{x}_i, \mathbf{x}_j \rangle = \mathbf{x}_i \cdot \mathbf{x}_j$$
- for example, to compute the output \mathbf{y}_3 in the above example, we need to compute $\text{score}(\mathbf{x}_3, \mathbf{x}_1)$, $\text{score}(\mathbf{x}_3, \mathbf{x}_2)$, $\text{score}(\mathbf{x}_3, \mathbf{x}_3)$

self-attention

- example (a simple self-attention):



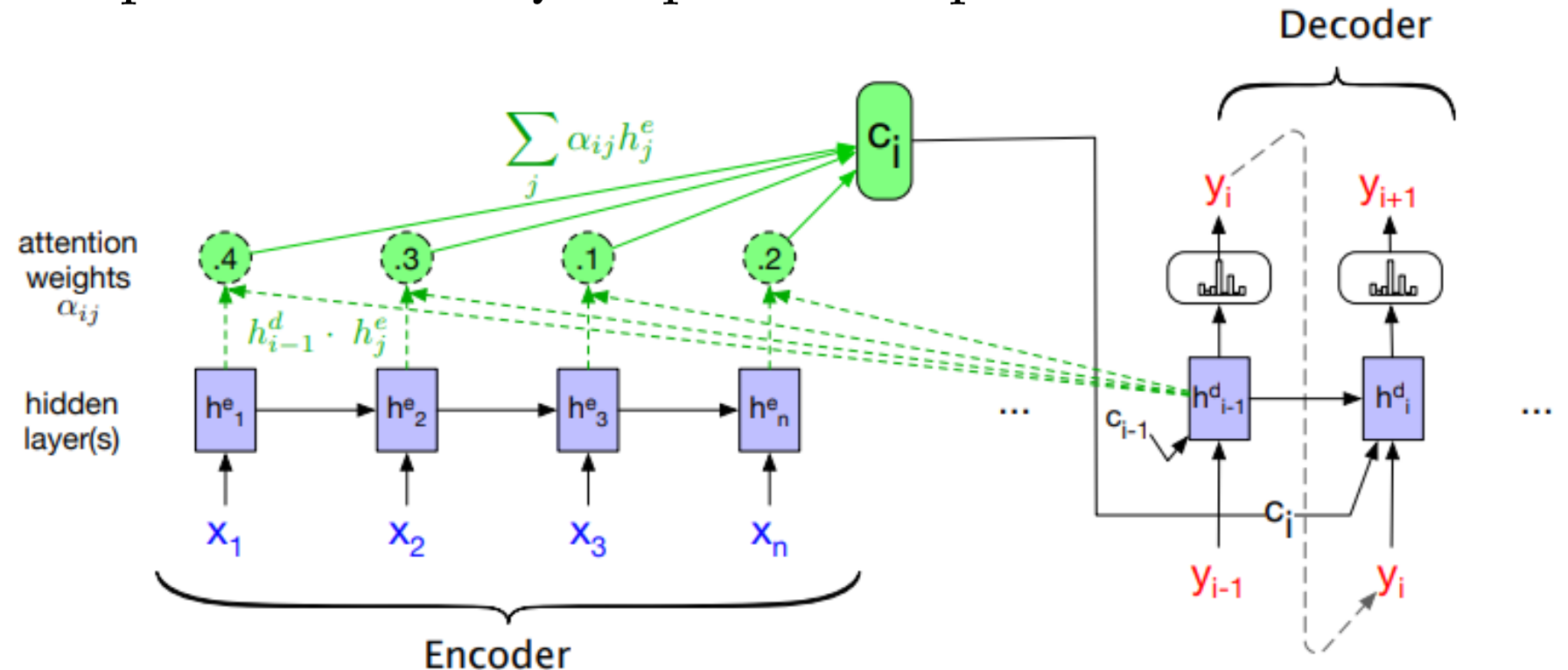
- the output y_3 then depends on a weighted average of x_1 , x_2 and x_3 ,
$$y_3 = \alpha_{31}x_1 + \alpha_{32}x_2 + \alpha_{33}x_3$$

- the weights α_{ij} are determined by these scores
- it is natural to normalize the scores using softmax, so that the weights are

$$\alpha_{ij} = \text{softmax}(\text{score}(x_i, x_j)) = \frac{\exp(\text{score}(x_i, x_j))}{\sum_{k=1}^i \exp(\text{score}(x_i, x_k))}, \quad \text{for all } j \leq i$$

self-attention

- an example of using the simple self-attention
- each output \mathbf{y}_i (more precisely, the hidden \mathbf{h}_i) in the decoder has a context \mathbf{c}_i which depends differently on previous input vectors



self-attention in transformer

- the self-attention regime in transformers is more complicated
- each input vector \mathbf{x}_i is transformed into three vectors of different roles:
 - \mathbf{q}_i (query): current input used in comparison
 - \mathbf{k}_i (key): previous input used in comparison
 - \mathbf{v}_i (value): input for taking weighted average
- they can be simply calculated by applying linear transformations:
 - $\mathbf{q}_i = \mathbf{W}_Q \mathbf{x}_i, \mathbf{k}_i = \mathbf{W}_K \mathbf{x}_i, \mathbf{v}_i = \mathbf{W}_V \mathbf{x}_i$

Freud's Structure of the Human Psyche



Id:

Instincts



Ego:

Reality



Superego:

Morality

self-attention in transformer

- if we make a simple analogy to the previous example, the scores are given by

$$\text{score}(\mathbf{x}_i, \mathbf{x}_j) = \langle \mathbf{q}_i, \mathbf{k}_j \rangle = \mathbf{q}_i \cdot \mathbf{k}_j, \quad i \leq j$$

- the output is then

$$\mathbf{y}_i = \sum_{j \leq i} \alpha_{ij} \mathbf{v}_j = \sum_{j \leq i} \text{softmax}(\mathbf{q}_i \cdot \mathbf{k}_j) \mathbf{v}_j$$

self-attention in transformer

- nevertheless, numerically it is often beneficial to normalize the dot product, so that

$$\text{score}(\mathbf{x}_i, \mathbf{x}_j) = \frac{\langle \mathbf{q}_i, \mathbf{k}_j \rangle}{\sqrt{d_k}} = \frac{\mathbf{q}_i \cdot \mathbf{k}_j}{\sqrt{d_k}}, \quad i \leq j$$

where d_k is the dimensionality of the key vector

self-attention in transformer

- numerically, it is more efficient to implement the above in matrix forms by combining the input vectors into a matrix $\mathbf{X} \in \mathbb{R}^{N \times d}$
- then we transform it using three matrices $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V \in \mathbb{R}^{d \times d}$
- we get three output matrices of size $N \times d$:

$$\mathbf{Q} = \mathbf{XW}_Q, \quad \mathbf{K} = \mathbf{XW}_K, \quad \mathbf{V} = \mathbf{XW}_V$$

self-attention in transformer

- in the matrix form, the weighted average of “values” can be written as

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

- however, we need to implement $\mathbf{Q}\mathbf{K}^T$ so that it is “causal” (that is, we calculate a score $\mathbf{q}_i \cdot \mathbf{k}_j$ only when $i \leq j$)
- we need to get rid of (mask out) the upper-triangular entries of $\mathbf{Q}\mathbf{K}^T$

self-attention in transformer

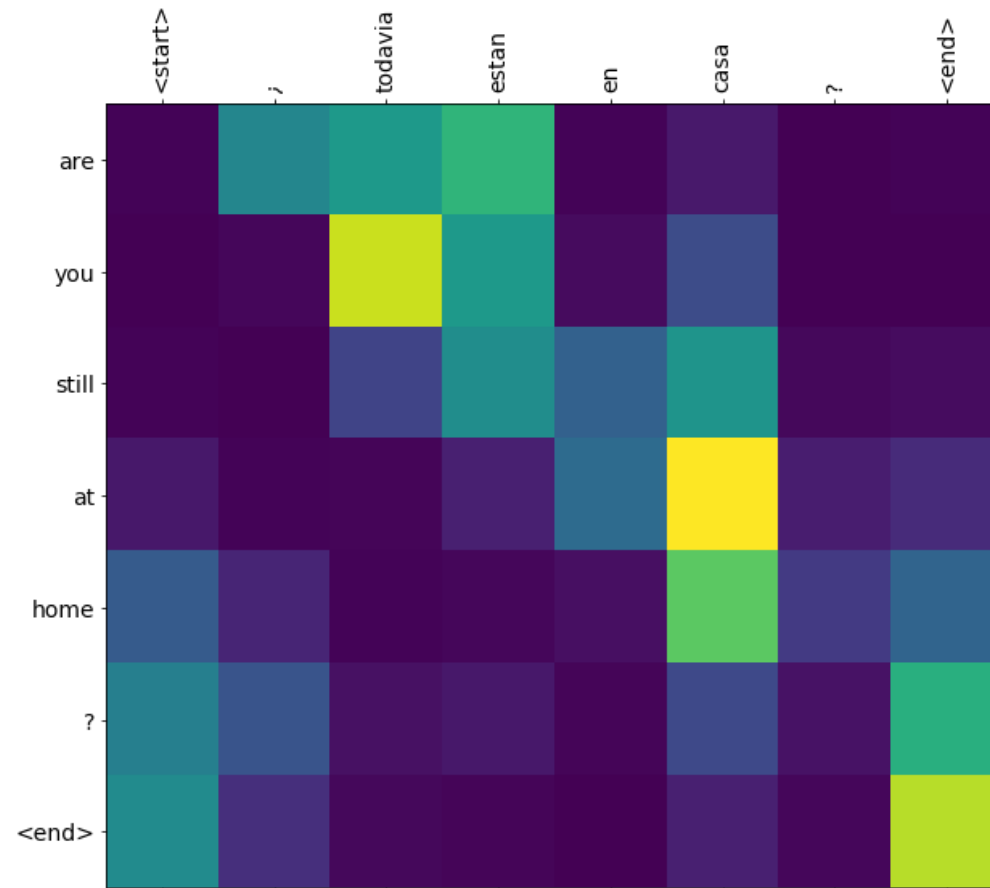
- we need to get rid of (mask out) the upper-triangular entries of \mathbf{QK}^T

N

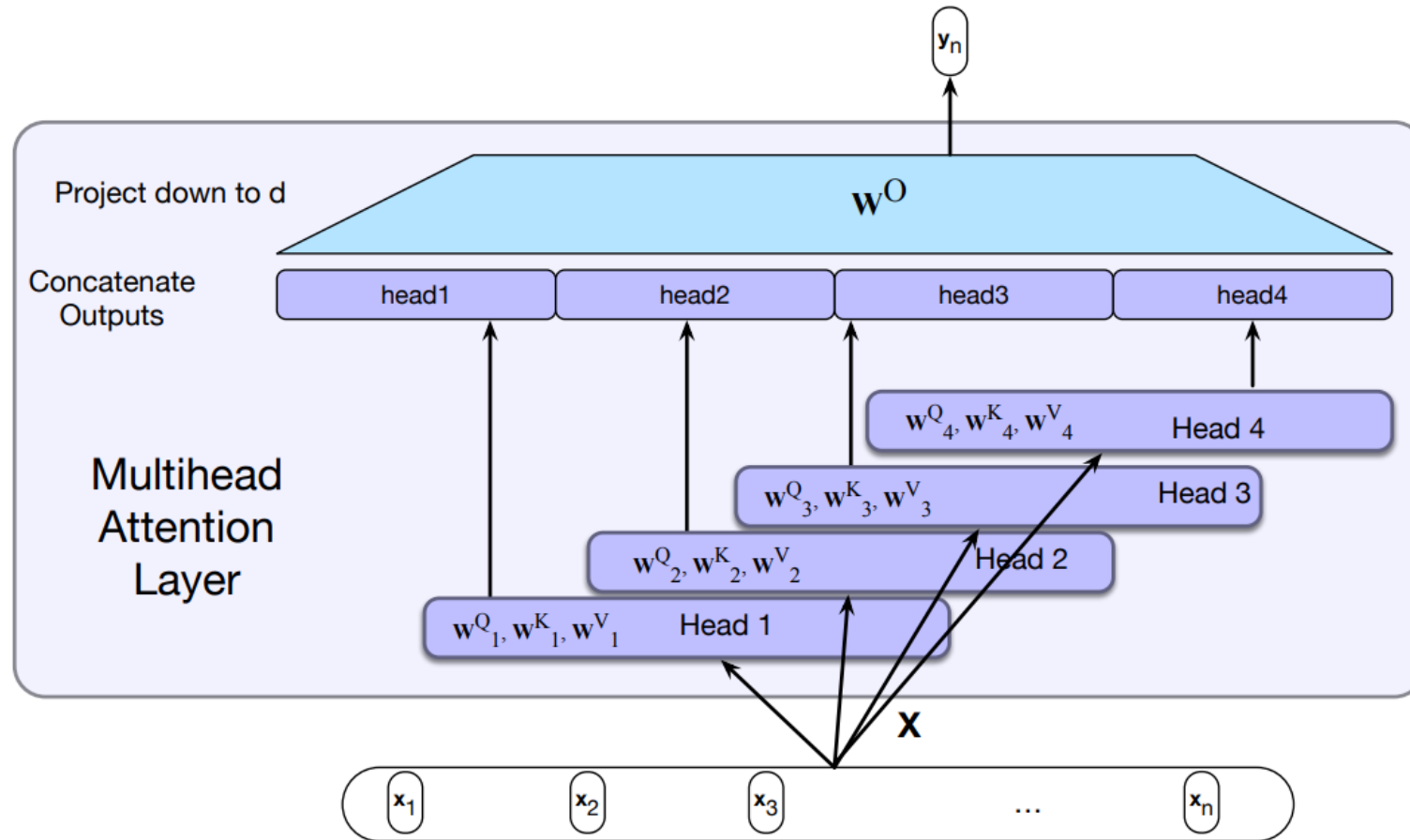
q1•k1	$-\infty$	$-\infty$	$-\infty$	$-\infty$
q2•k1	q2•k2	$-\infty$	$-\infty$	$-\infty$
q3•k1	q3•k2	q3•k3	$-\infty$	$-\infty$
q4•k1	q4•k2	q4•k3	q4•k4	$-\infty$
q5•k1	q5•k2	q5•k3	q5•k4	q5•k5

N

an example of attention weights

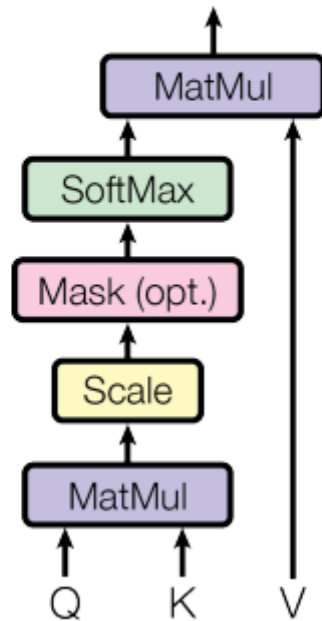


multi-head attention



multi-head attention

Scaled Dot-Product Attention



Multi-Head Attention

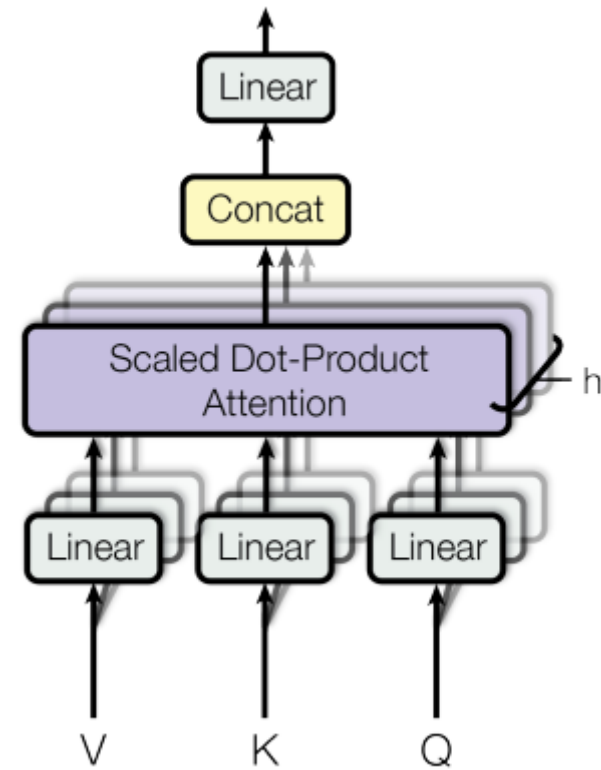
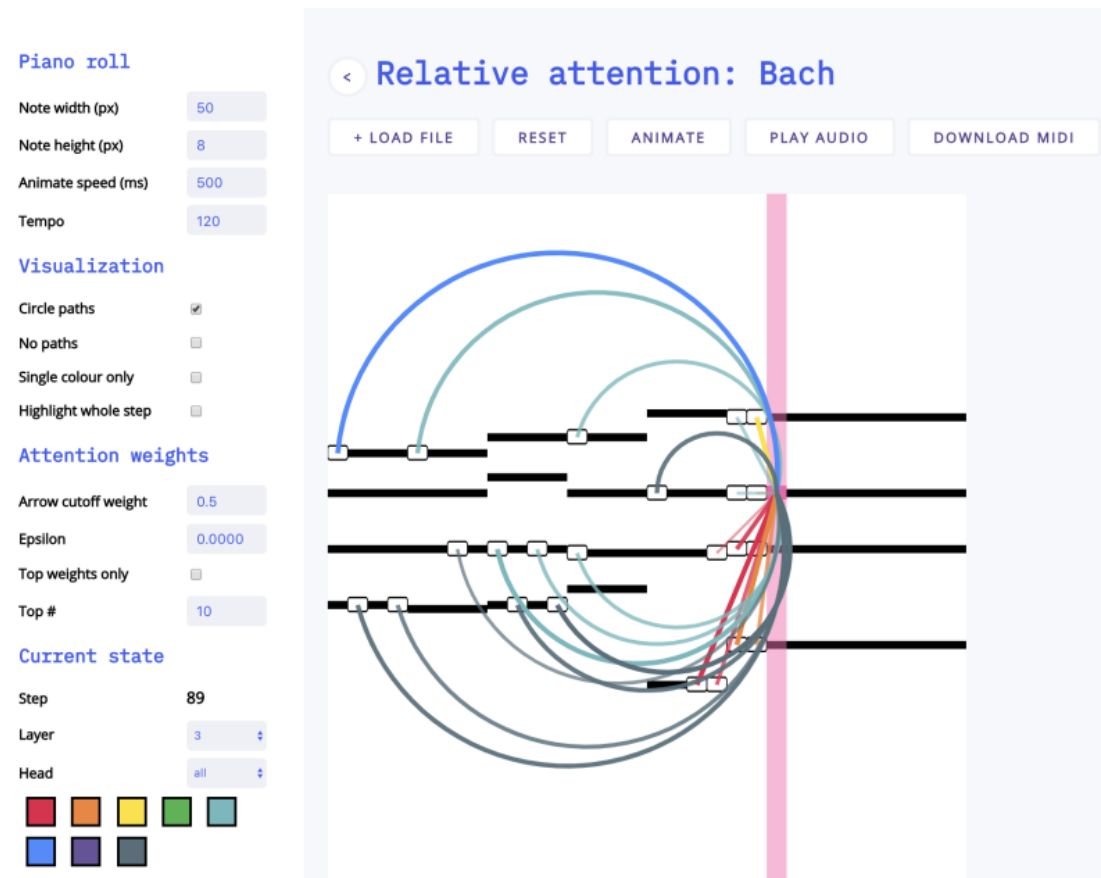


image credit: Vaswani, A. et al (2017). Attention is all you need.

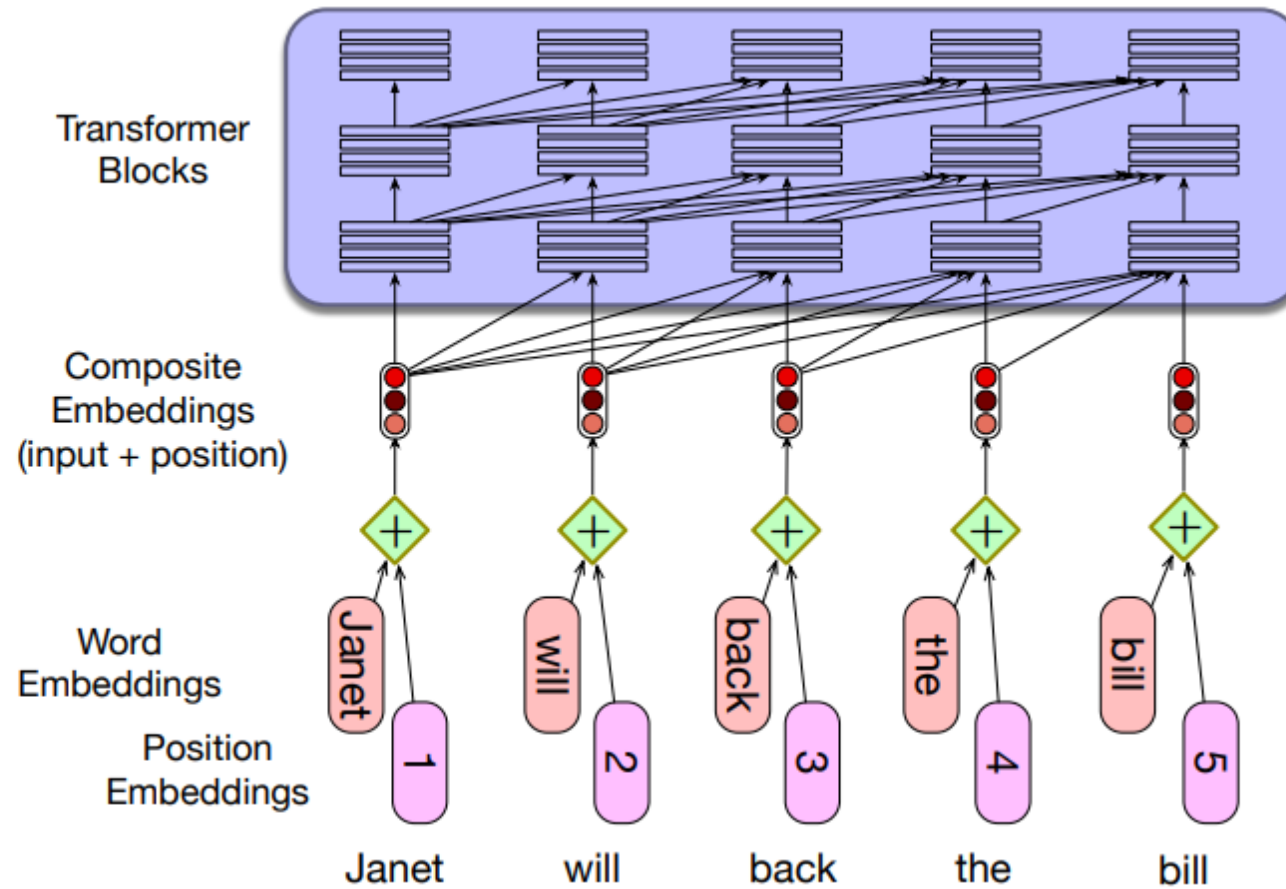
another example

- <https://storage.googleapis.com/nips-workshop-visualization/index.html>



7.2 transformer

positional embedding (encoding)



positional embedding (encoding)

- transformers process input data in parallel, which means they don't inherently understand the order or position of elements in a sequence
- positional embedding refers to the technique used to inject information about the position of words or tokens in a sequence into the model
- create a matrix where each row corresponds to a position in the sequence; these embeddings are then **added** element-wise to the input embeddings

positional embedding (encoding)

- one popular choice

$$\vec{p}_t^{(i)} = f(t)^{(i)} := \begin{cases} \sin(\omega_k \cdot t), & \text{if } i = 2k \\ \cos(\omega_k \cdot t), & \text{if } i = 2k + 1 \end{cases}$$

position

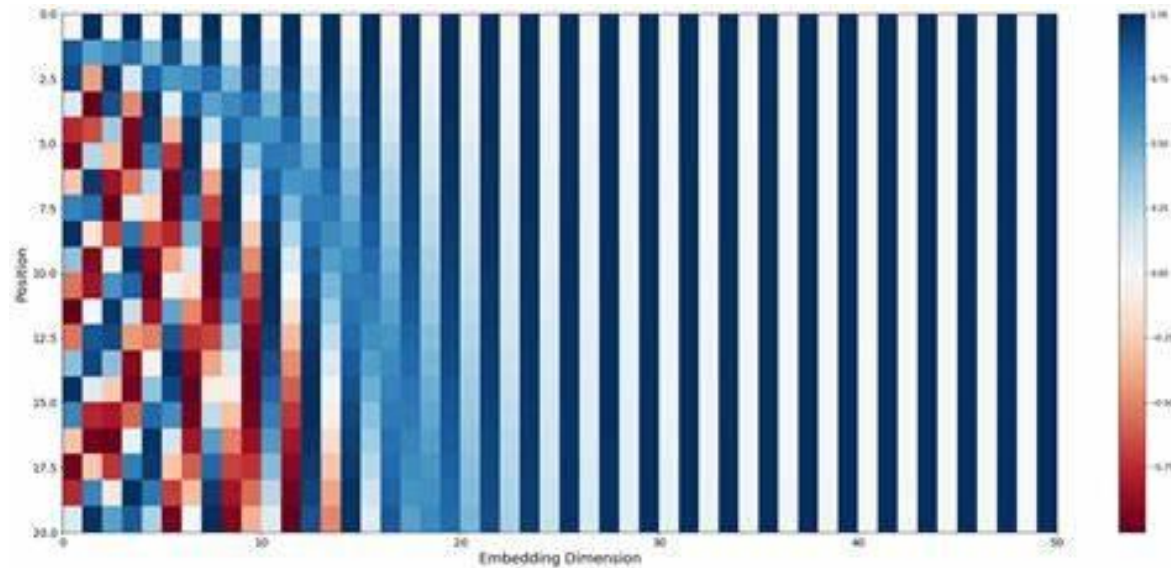
dimension

$$\omega_k = \frac{1}{10000^{2k/d}}$$

positional embedding (encoding)

- one popular choice

$$\vec{p}_t^{(i)} = f(t)^{(i)} := \begin{cases} \sin(\omega_k \cdot t), & \text{if } i = 2k \\ \cos(\omega_k \cdot t), & \text{if } i = 2k + 1 \end{cases}$$



positional embedding (encoding)

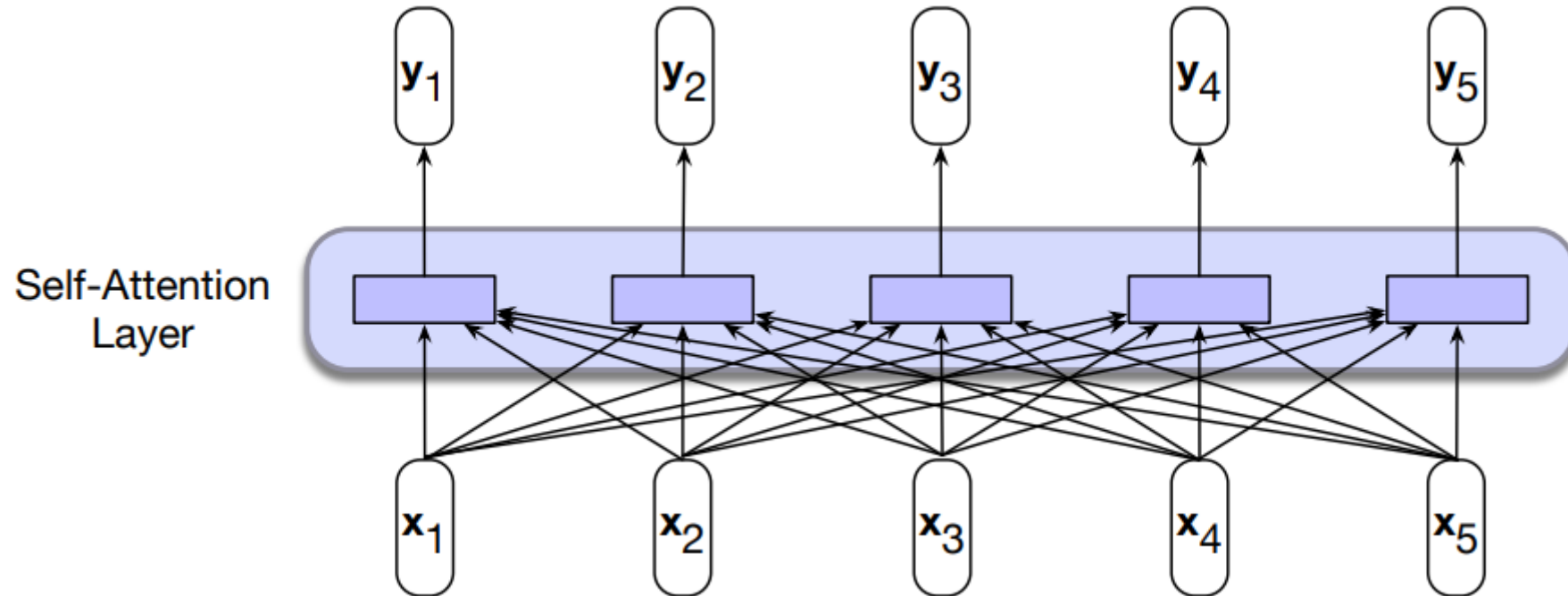
- it is also possible to encode relative positions of pairs of input
- it is also possible to encode positions for query and key

$$\text{RoPE}(q)[2i] = q[2i] \cdot \cos(\theta_i) - q[2i + 1] \cdot \sin(\theta_i)$$

$$\text{RoPE}(q)[2i + 1] = q[2i] \cdot \sin(\theta_i) + q[2i + 1] \cdot \cos(\theta_i)$$

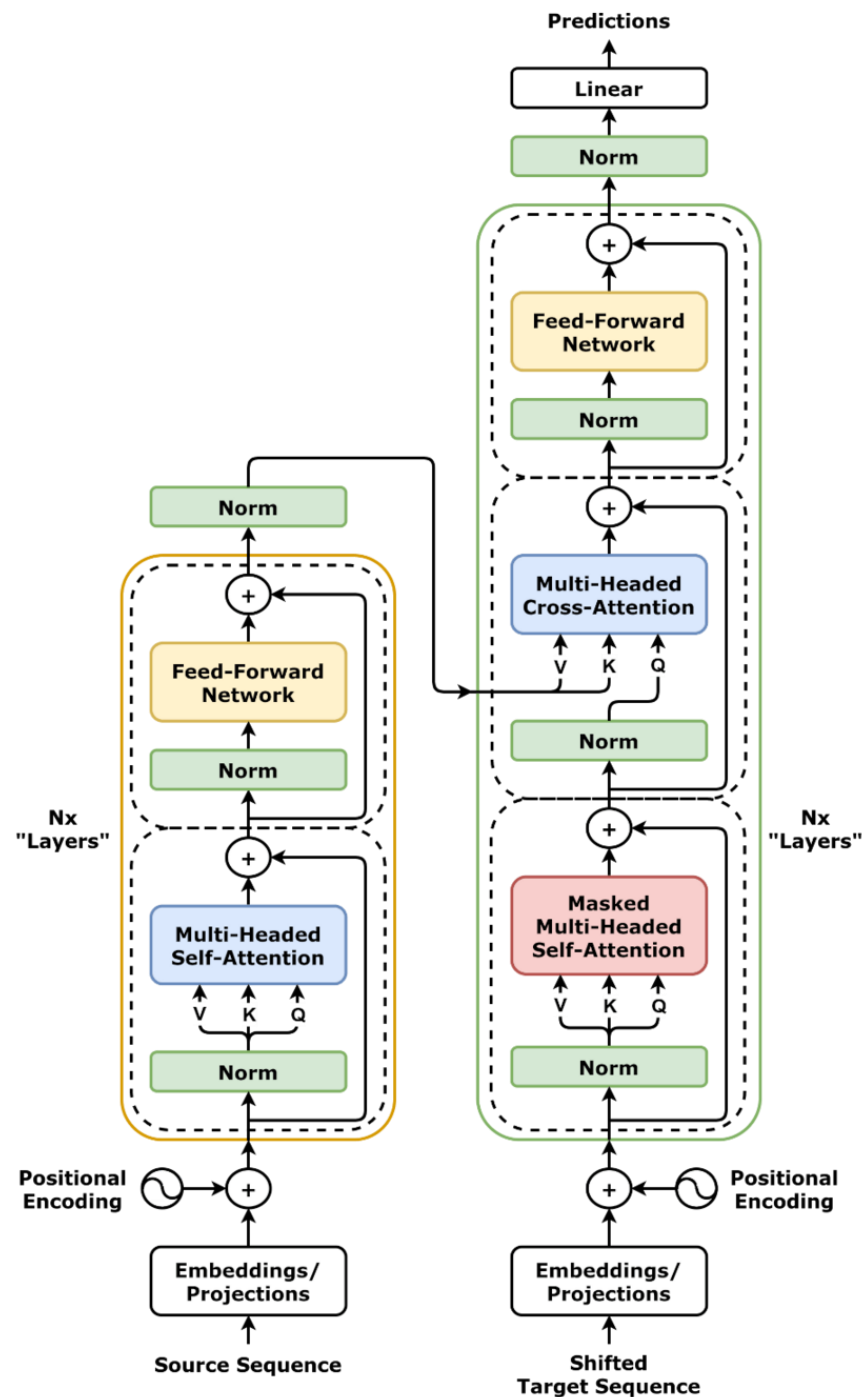
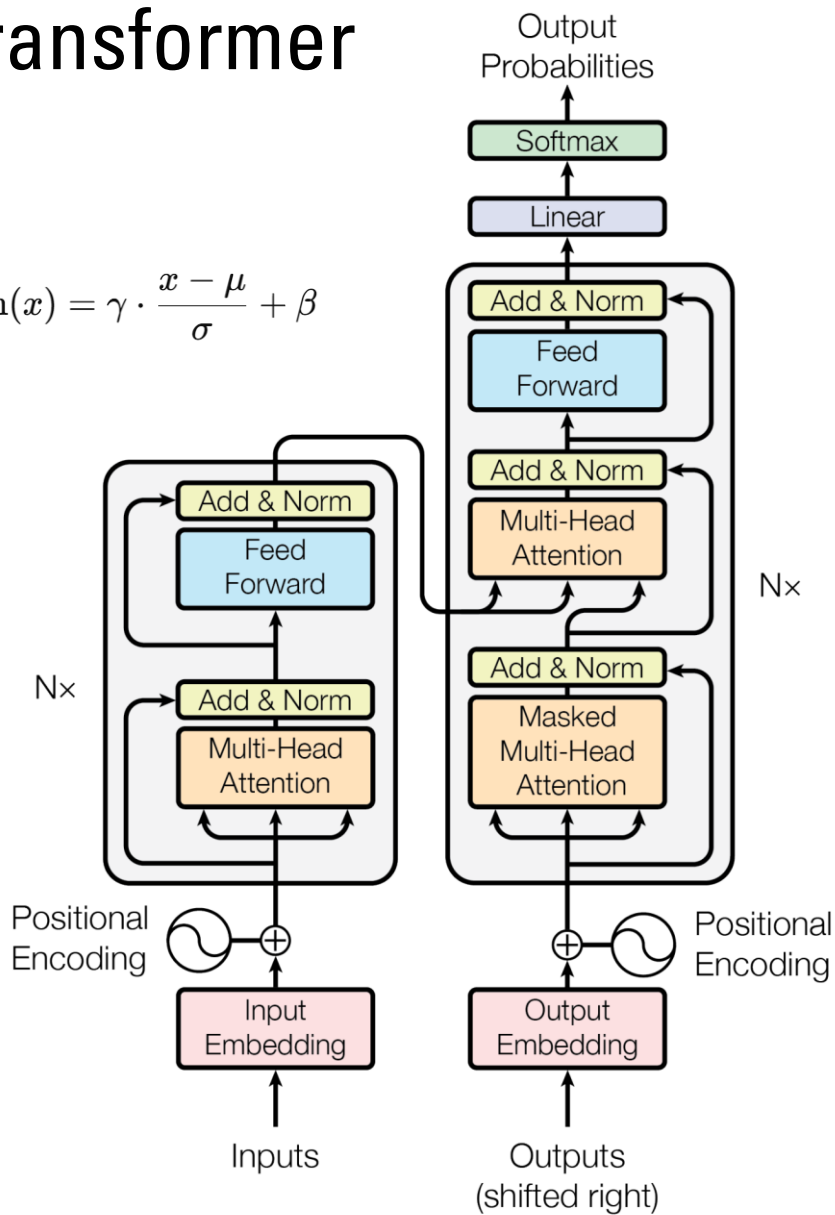
$$\theta_i = p/10000^{2i/d} \quad (\text{like classic sinusoidal freq})$$

non-causal attention is also possible



transformer

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

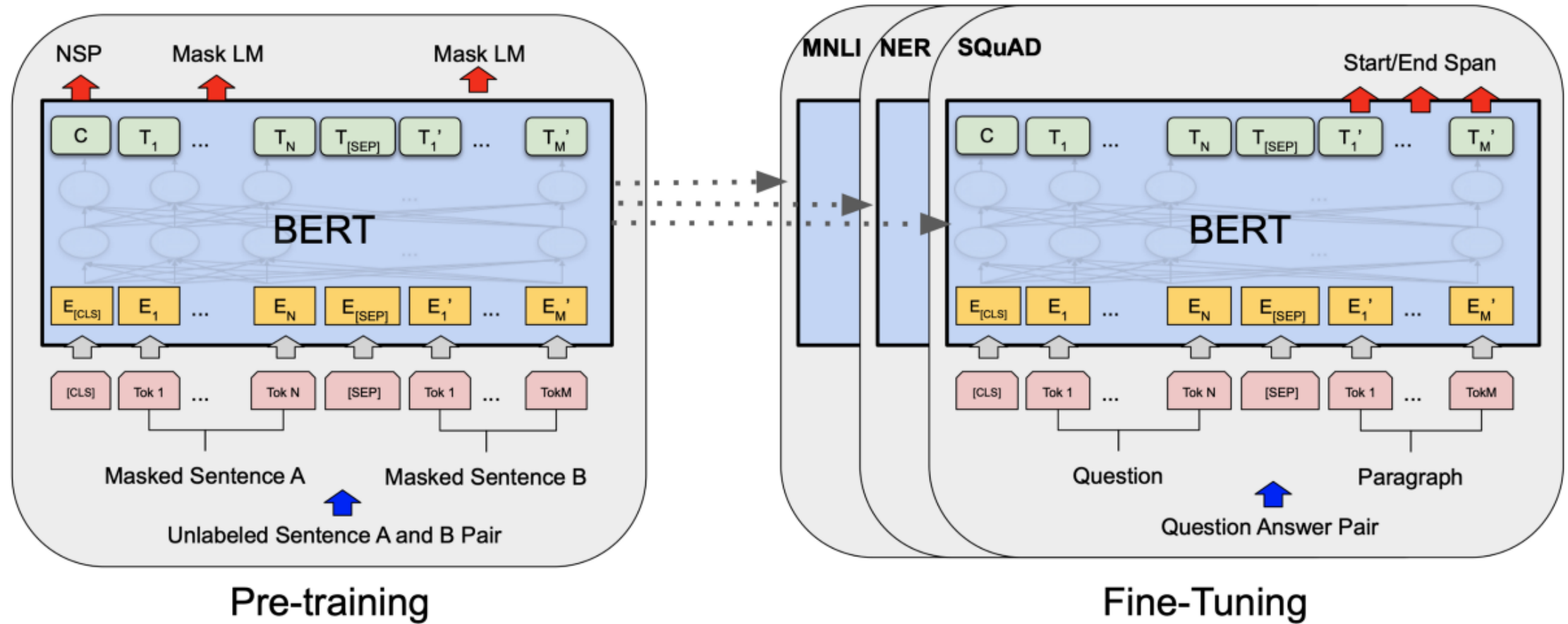


- [AnnotatedTransformer.ipynb - Colab](#)

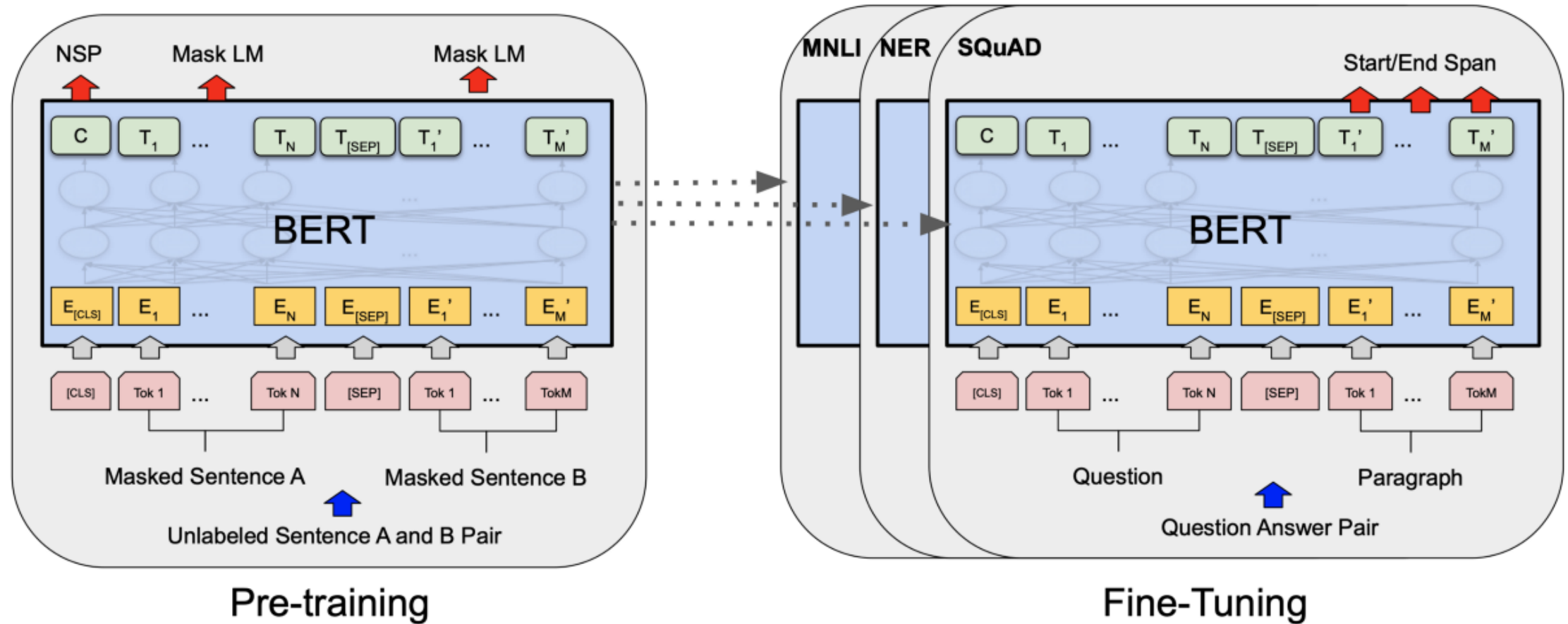
advantage of transformers

- parallelization and efficiency
- long-range dependencies
- reduced vanishing gradient problem
- interpretability
- ...

example: BERT (Bidirectional Encoder Repr. fr. Transformers)



example: BERT (Bidirectional Encoder Repr. fr. Transformers)



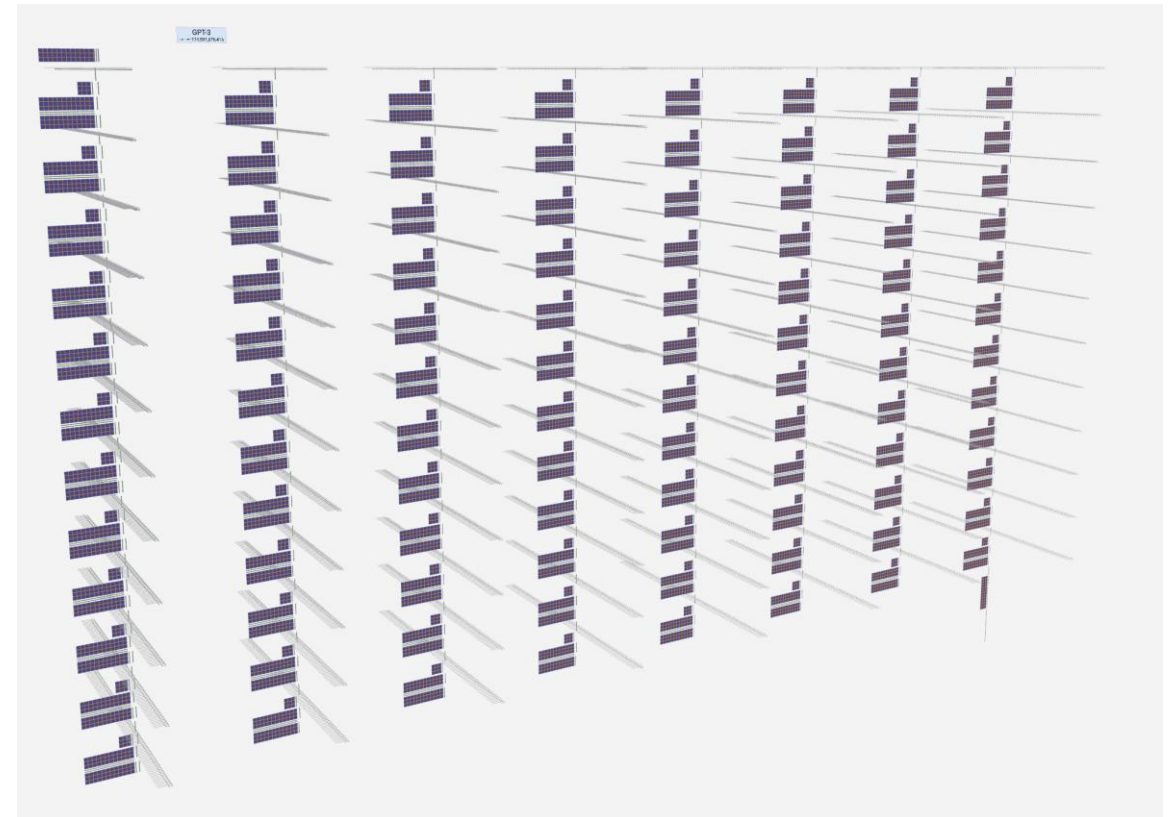
It uses only the encoder part of the transformer.

Not "causal": bi-directional.

example: GPT (Generative Pretrained Transformer)

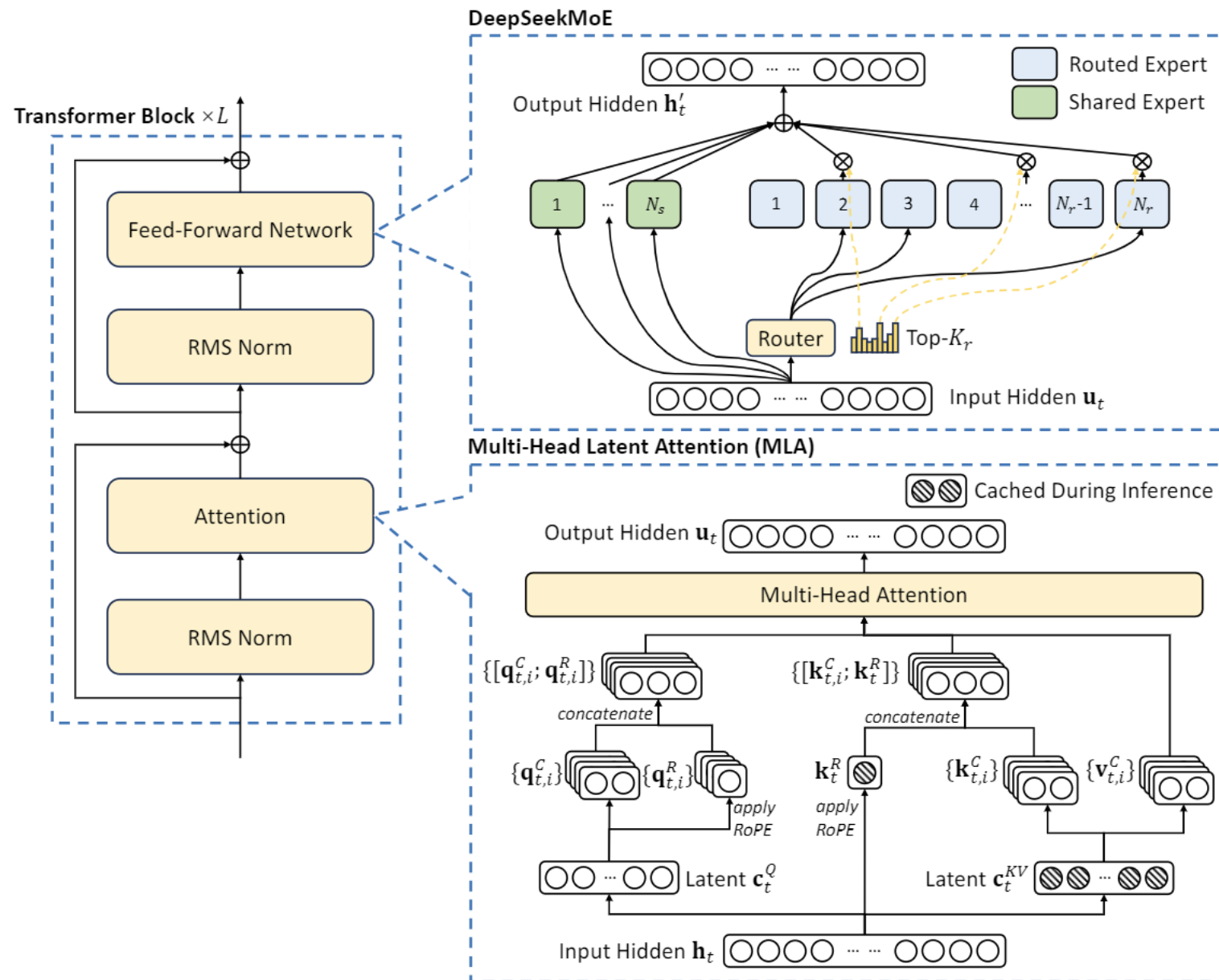
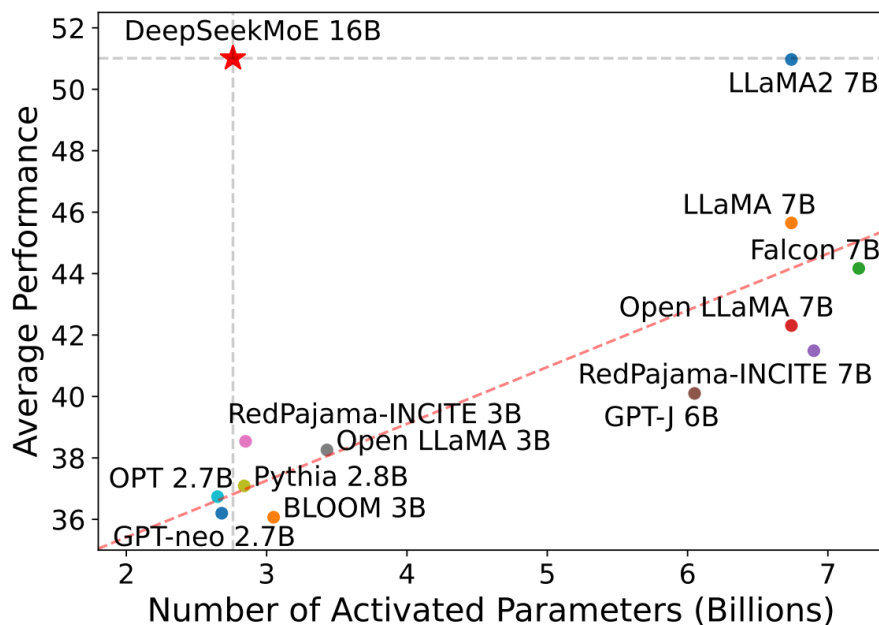
- decoder-only

Model Name	n_{params}	n_{layers}	d_{model}	n_{heads}	d_{head}	Batch Size	Learning Rate
GPT-3 Small	125M	12	768	12	64	0.5M	6.0×10^{-4}
GPT-3 Medium	350M	24	1024	16	64	0.5M	3.0×10^{-4}
GPT-3 Large	760M	24	1536	16	96	0.5M	2.5×10^{-4}
GPT-3 XL	1.3B	24	2048	24	128	1M	2.0×10^{-4}
GPT-3 2.7B	2.7B	32	2560	32	80	1M	1.6×10^{-4}
GPT-3 6.7B	6.7B	32	4096	32	128	2M	1.2×10^{-4}
GPT-3 13B	13.0B	40	5140	40	128	2M	1.0×10^{-4}
GPT-3 175B or “GPT-3”	175.0B	96	12288	96	128	3.2M	0.6×10^{-4}



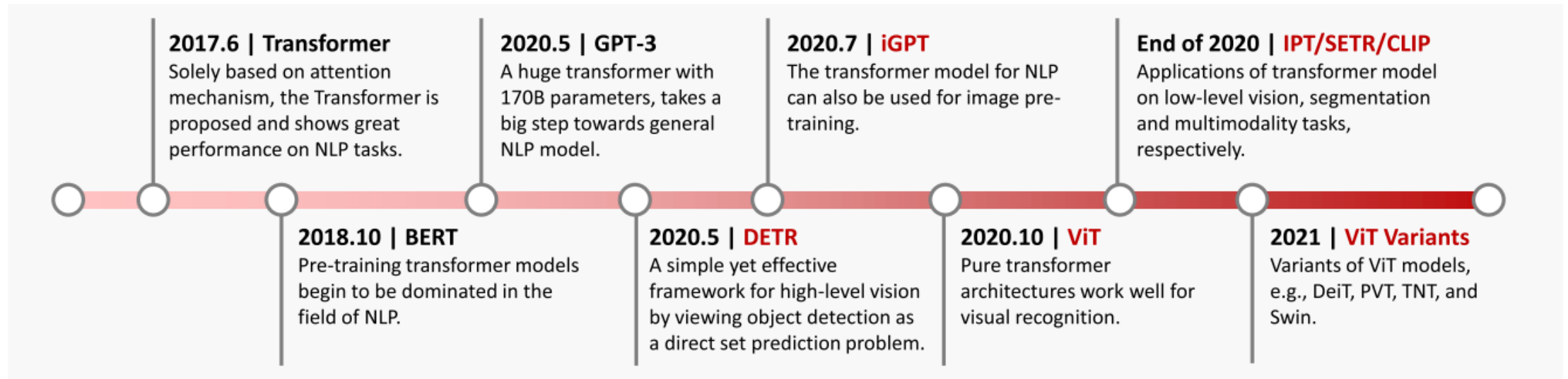
example: DeepSeek

- decoder-only

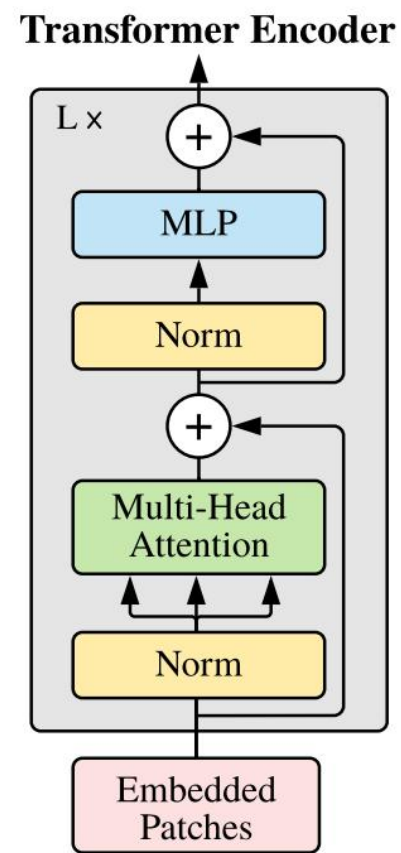
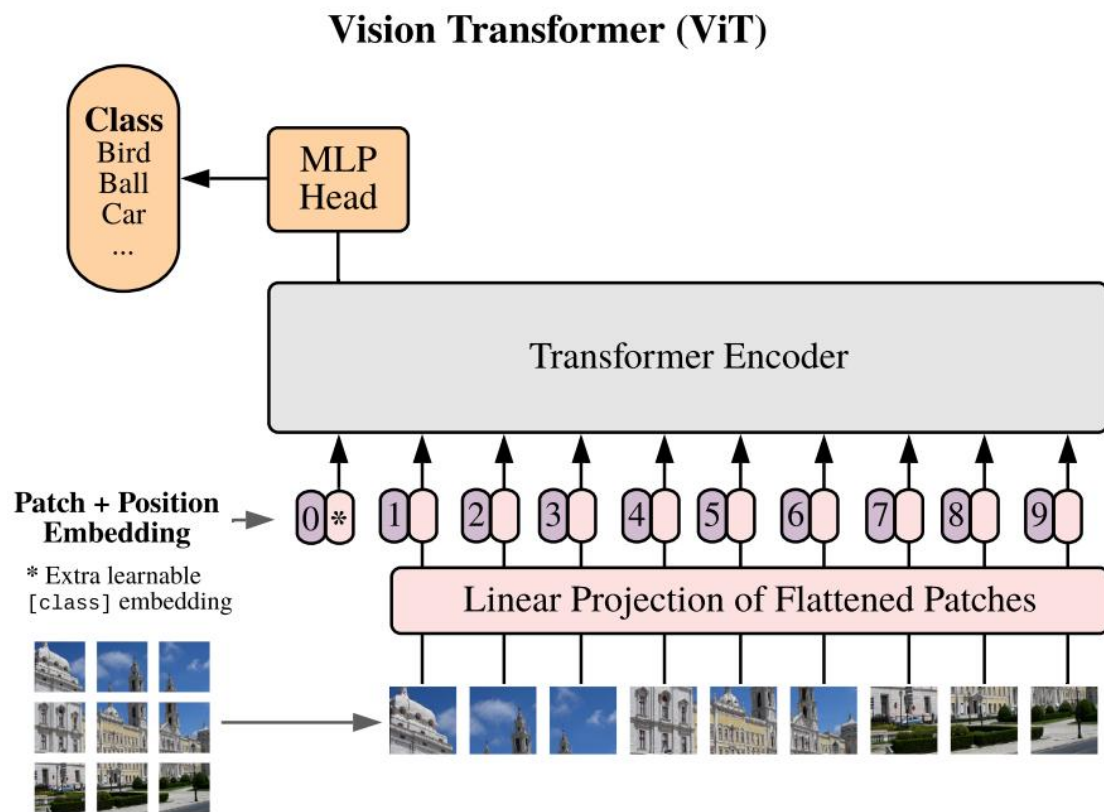


7.3 vision transformer

a history



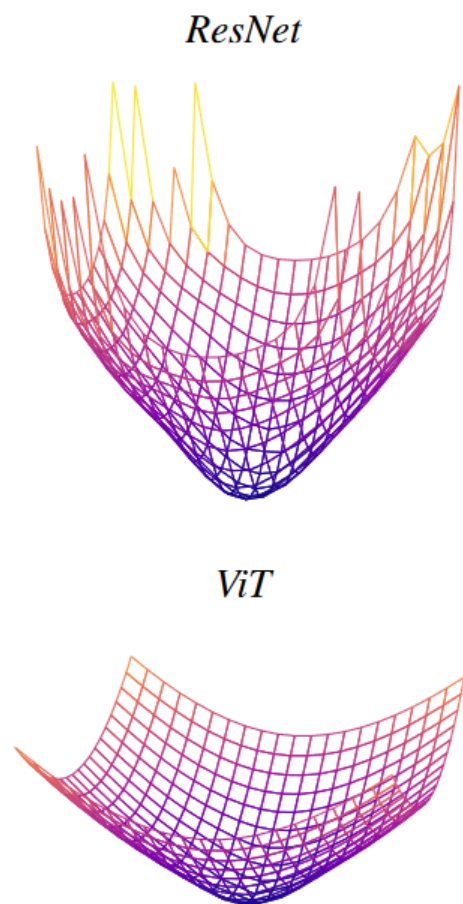
vision transformer



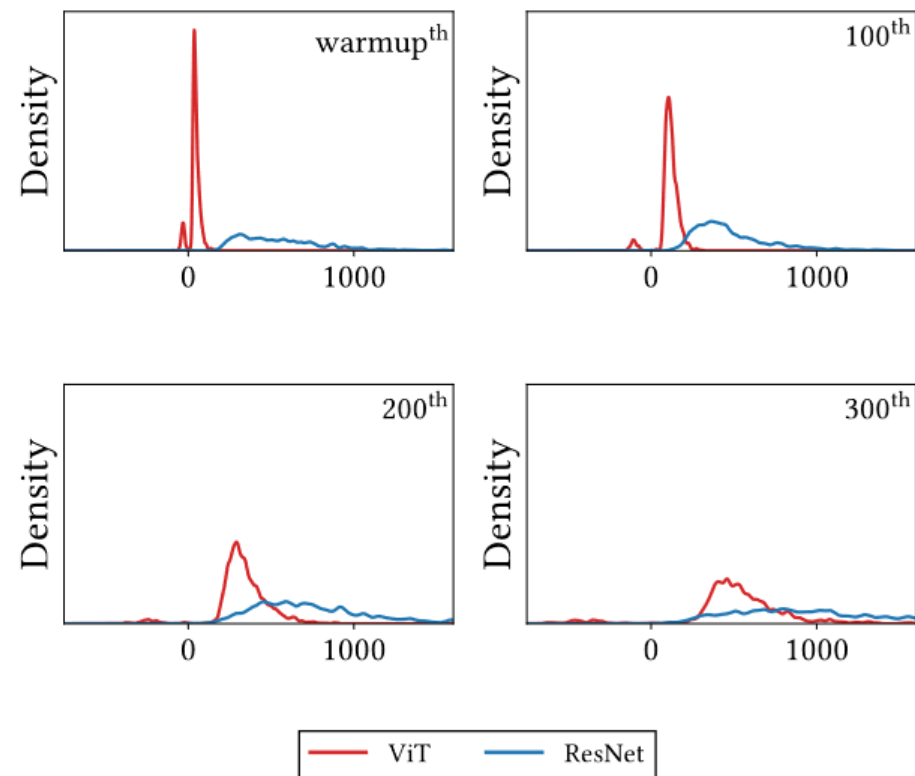
advantages of vision transformers

- global context understanding
- scalability
- parameter efficiency
- ...

How does vision transformer work?

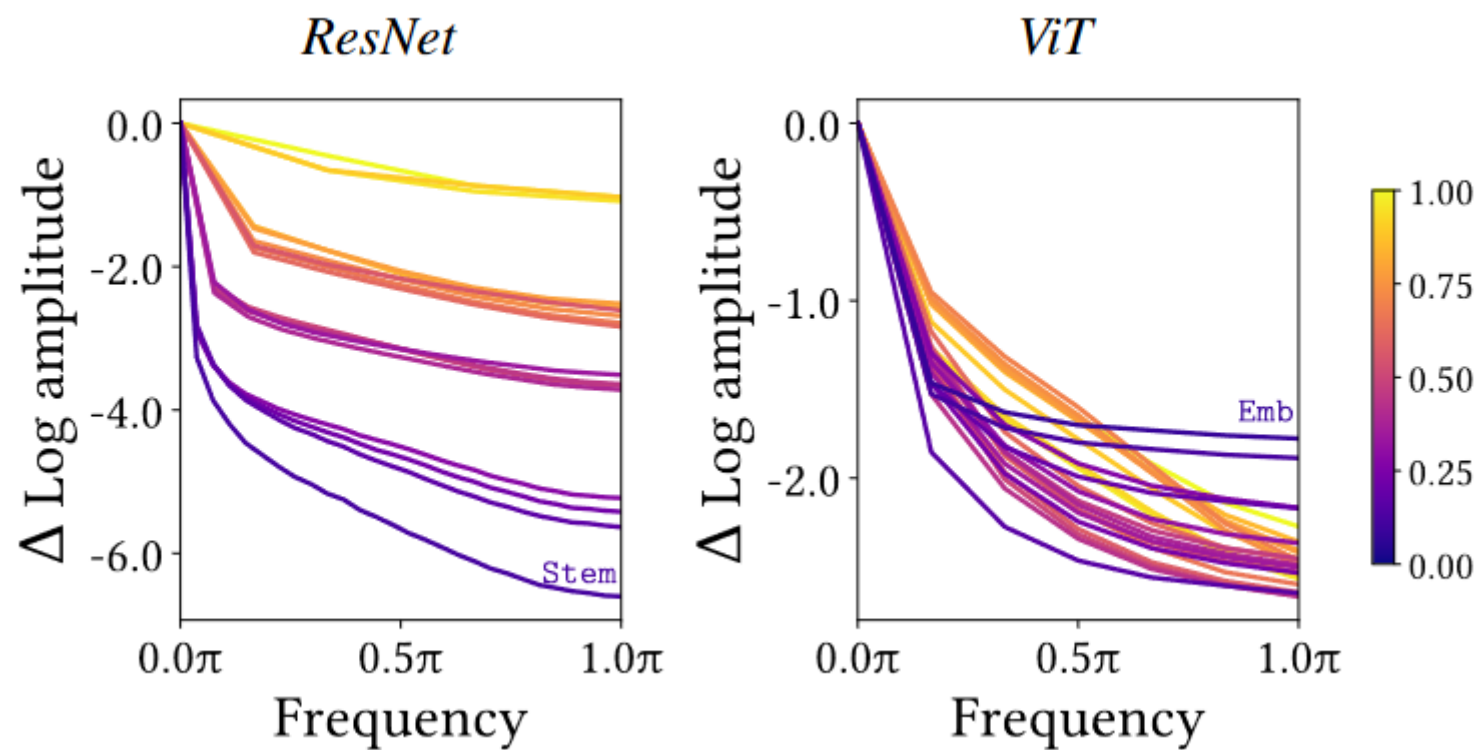


(a) Loss landscape visualizations

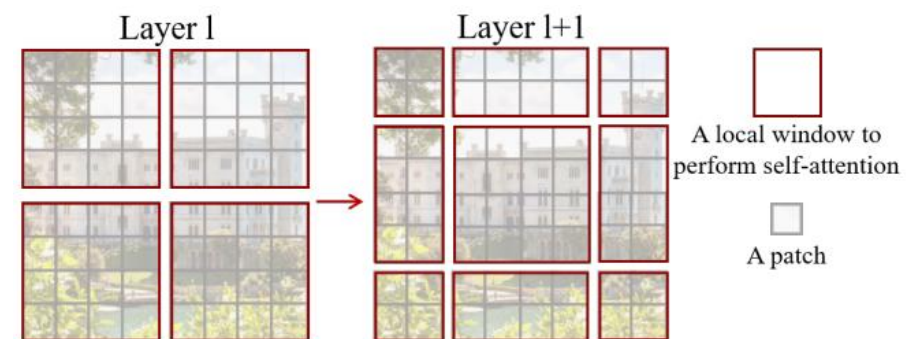
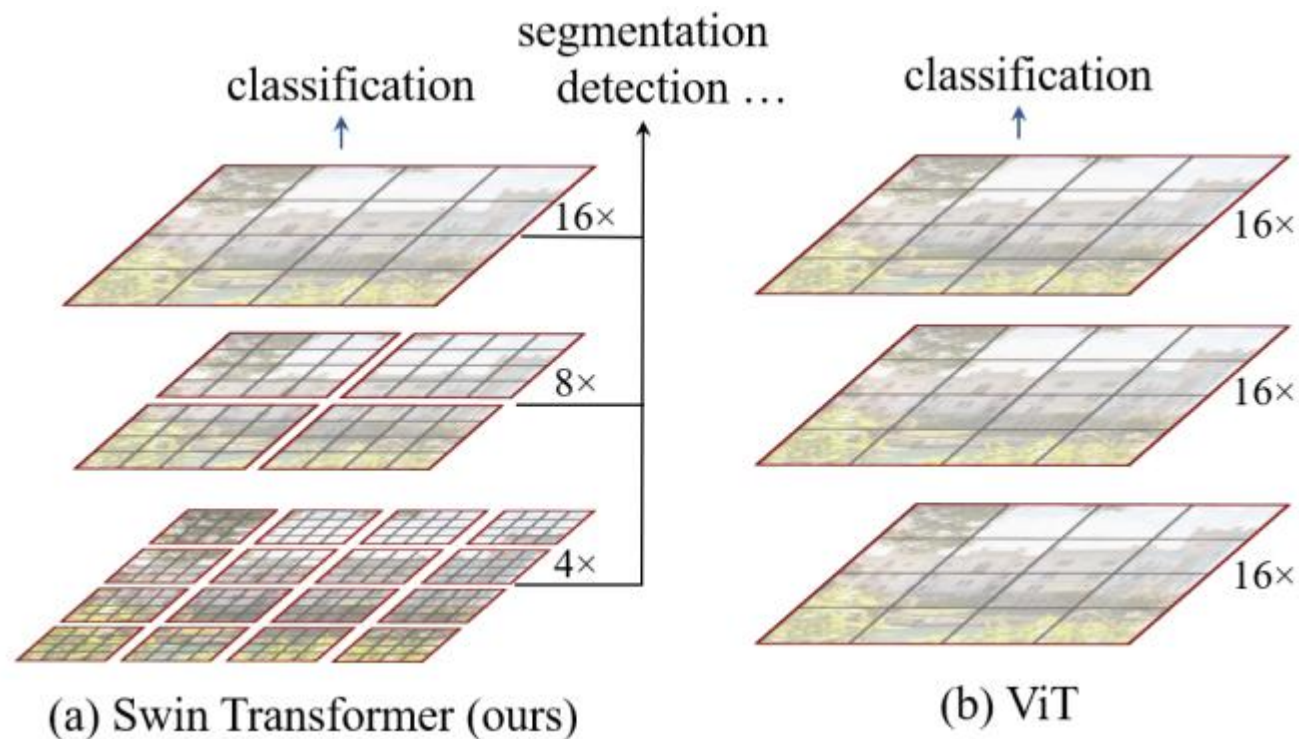


(b) Hessian max eigenvalue spectra

How does vision transformer work?



swin (shifted windows) transformer



swin (shifted windows) transformer

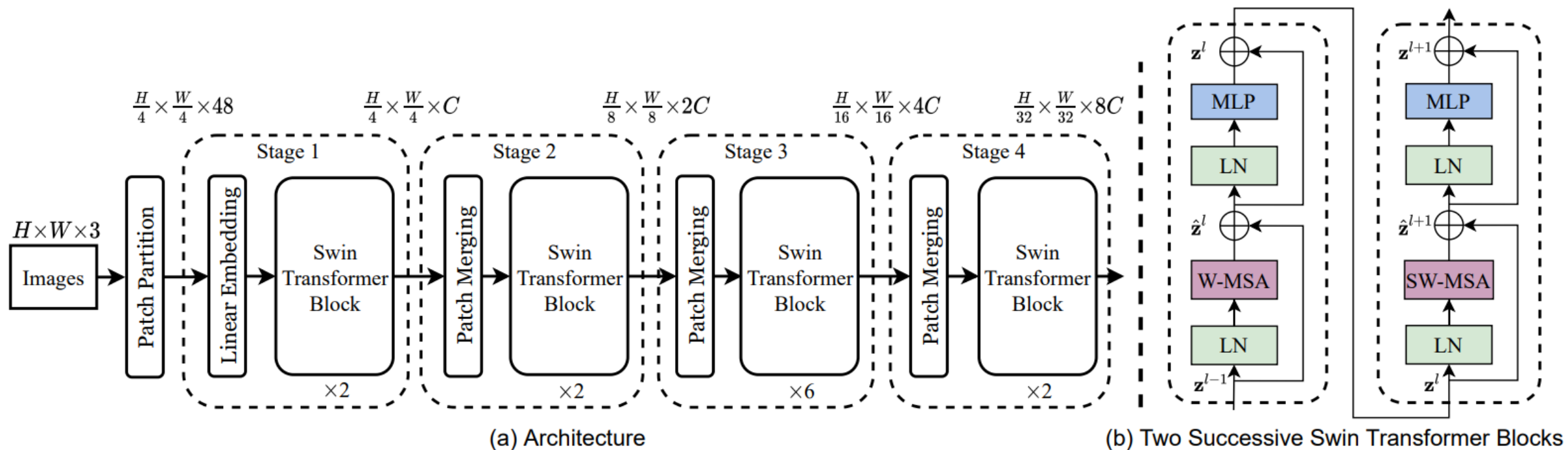


Figure 3. (a) The architecture of a Swin Transformer (Swin-T); (b) two successive Swin Transformer Blocks (notation presented with Eq. (3)). W-MSA and SW-MSA are multi-head self attention modules with regular and shifted windowing configurations, respectively.

Thank you!

Reference

- Ch 10, 11 & 16, Jurafsky, D., & Martin, J. H. (2023). *Speech and language processing (draft)*. <https://web.stanford.edu/~jurafsky/slp3/>
- Goldberg, Y. (2016). A primer on neural network models for natural language processing. *Journal of Artificial Intelligence Research*, 57, 345-420.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. *Advances in neural information processing systems*, 30.
- Huang, C. Z. A., Dinculescu, M., Vaswani, A., & Eck, D. (2018). Visualizing music transformer. In *NeurIPS Workshop on Interpretability and Robustness in Audio, Speech, and Language*.
- Kenton, J. D. M. W. C., & Toutanova, L. K. (2019). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of NAACL-HLT* (pp. 4171-4186).
- Han, K., Wang, Y., Chen, H., Chen, X., Guo, J., Liu, Z., ... & Tao, D. (2022). A survey on vision transformer. *IEEE transactions on pattern analysis and machine intelligence*, 45(1), 87-110.
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., ... & Houlsby, N. (2020). An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR 2021*.
- Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., ... & Guo, B. (2021). Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision* (pp. 10012-10022).
- Cao, H., Wang, Y., Chen, J., Jiang, D., Zhang, X., Tian, Q., & Wang, M. (2022). Swin-UNet: Unet-like pure transformer for medical image segmentation. In *ECCV 2022*.

