

Transformers

Ju Sun

Computer Science & Engineering

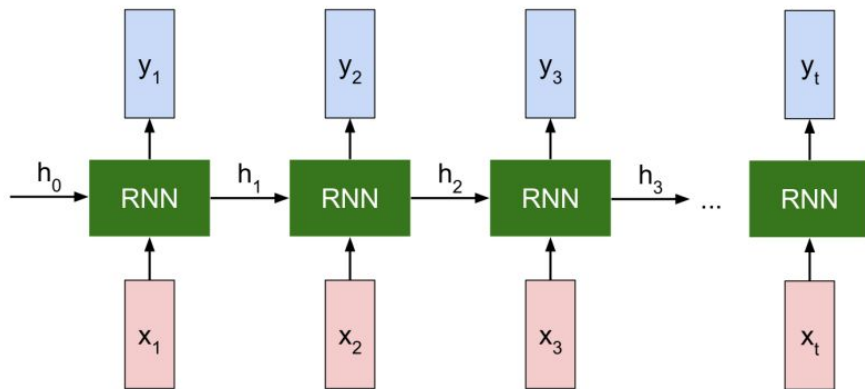
Nov 22, 2022



UNIVERSITY OF MINNESOTA
Driven to DiscoverSM

Recap

RNN: model sequences



(Credit: Stanford CS231N)

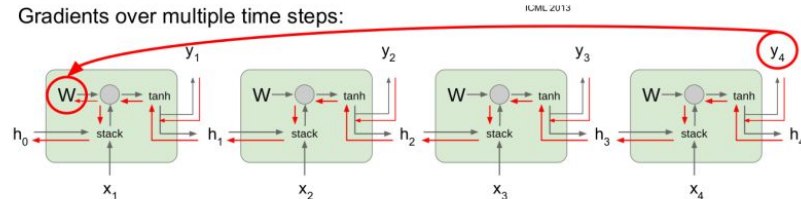
$$\mathbf{h}_t = \tanh(\mathbf{W}_h \mathbf{h}_{t-1} + \mathbf{W}_x \mathbf{x}_t)$$

$$\mathbf{y}_t = \mathbf{V}_y \mathbf{h}_t$$

\mathbf{W}_h , \mathbf{W}_x and \mathbf{V}_y are shared across the sequence

Vanishing/exploding gradient issue

Gradients over multiple time steps:



(Credit: Stanford CS231N)

$$\begin{aligned} \frac{\partial L_t}{\partial \mathbf{W}} &= \frac{\partial L_t}{\partial \mathbf{h}_t} \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} \dots \frac{\partial \mathbf{h}_1}{\partial \mathbf{W}} = \frac{\partial L_t}{\partial \mathbf{h}_t} \left(\prod_{k=2}^t \frac{\partial \mathbf{h}_k}{\partial \mathbf{h}_{k-1}} \right) \frac{\partial \mathbf{h}_1}{\partial \mathbf{W}} \\ &= \frac{\partial L_t}{\partial \mathbf{h}_t} \left(\prod_{k=2}^t \text{diag}(\tanh'(\mathbf{W}_h \mathbf{h}_{k-1} + \mathbf{W}_x \mathbf{x}_k)) \mathbf{W}_h \right) \frac{\partial \mathbf{h}_1}{\partial \mathbf{W}} \end{aligned}$$

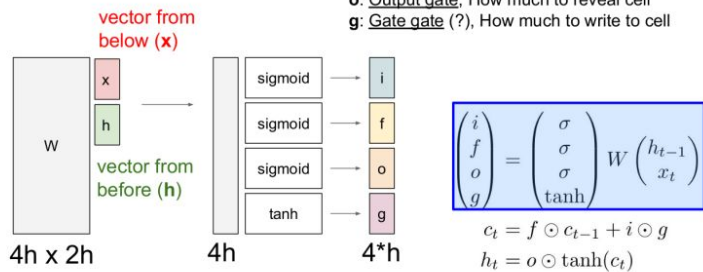
- * when $\|\mathbf{W}_h\| > 1$, gradient **explodes** if t large
- * when $\|\mathbf{W}_h\| < 1$, gradient **vanishes** if t large

$$\begin{aligned} &\left\| \prod_{k=2}^t \text{diag}(\tanh'(\mathbf{W}_h \mathbf{h}_{k-1} + \mathbf{W}_x \mathbf{x}_k)) \mathbf{W}_h \right\| \\ &\leq \prod_{k=2}^t \left\| \text{diag}(\tanh'(\mathbf{W}_h \mathbf{h}_{k-1} + \mathbf{W}_x \mathbf{x}_k)) \right\| \|\mathbf{W}_h\| \\ &\leq \prod_{k=2}^t \left\| \text{diag}(\tanh'(\mathbf{W}_h \mathbf{h}_{k-1} + \mathbf{W}_x \mathbf{x}_k)) \right\| \|\mathbf{W}_h\|^{t-1} \end{aligned} \quad 2$$

Gated RNNs

Long Short Term Memory (LSTM)

[Hochreiter et al., 1997]



(Credit: Stanford CS231N)

u : **update gate**, control state update

r : **reset gate**, control how previous state affects new content

g : new content

Gated recurrent unit (GRU)

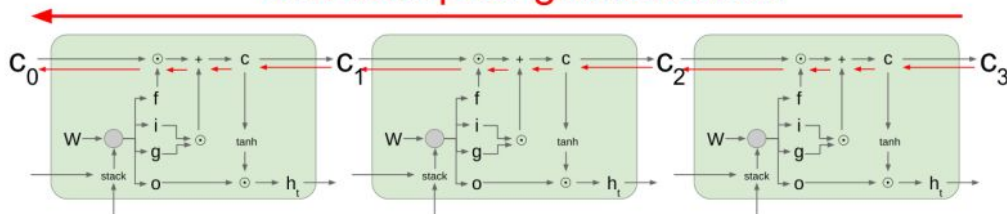
$$\begin{bmatrix} u \\ r \\ g \end{bmatrix} = \begin{bmatrix} \sigma \\ \sigma \end{bmatrix} \left(W \begin{bmatrix} h_{t-1} \\ x_t \end{bmatrix} \right)$$

$$g = \tanh(W_h (r \odot h_{t-1}) + W_x x_t + b_g)$$

$$h_t = u \odot h_{t-1} + (1 - u) \odot g$$

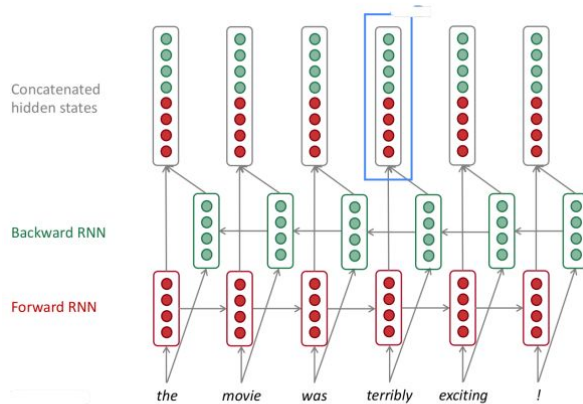
f, i, o are merged

Uninterrupted gradient flow!



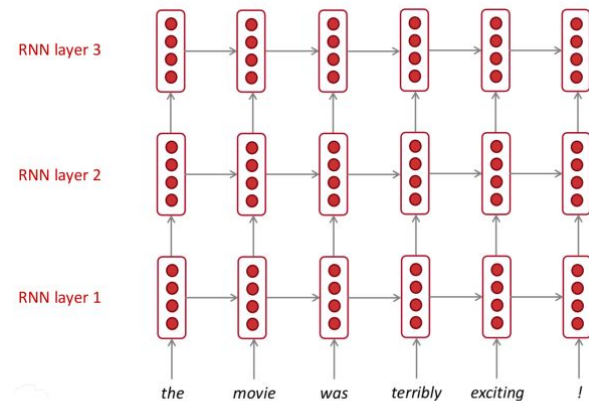
(Credit: Stanford CS231N)

Modern RNNs



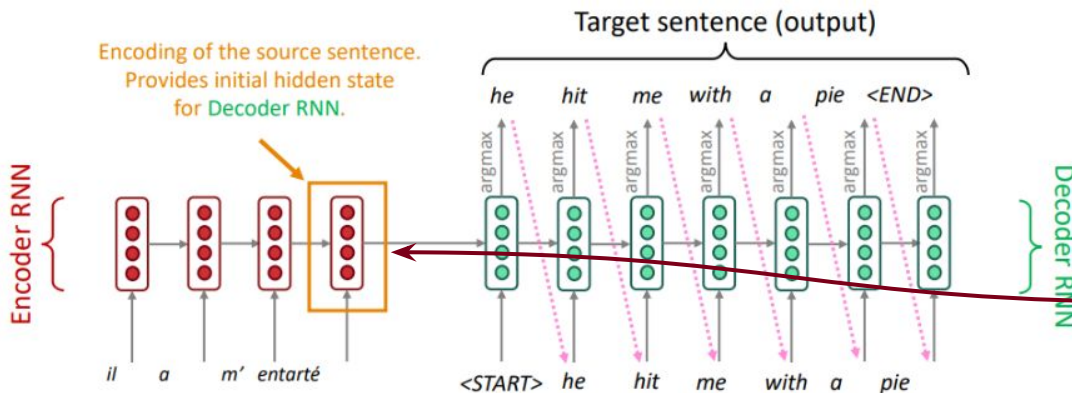
(Credit: Stanford CS224N)

Bidirectional RNN



(Credit: Stanford CS231N)

Deep RNN

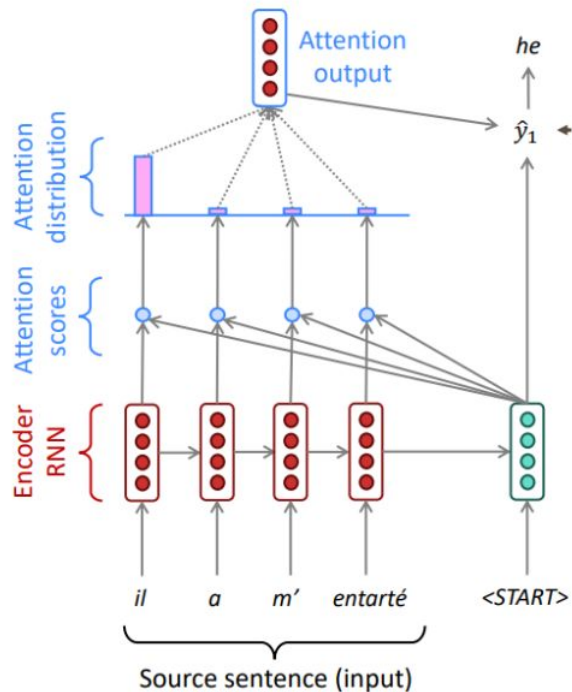


(Credit: Stanford CS231N)

Seq2Seq model

Bottleneck problem

Attention mechanism



(Credit: Stanford CS231N)

Input: source vectors $s_1, \dots, s_N \in \mathbb{R}^h$, and target vector t

Output: weighted summation

$$\sum_{j=1}^N w_j s_j \quad \text{where } w_j = \text{similarity}(s_j, t)$$

Many possibilities:

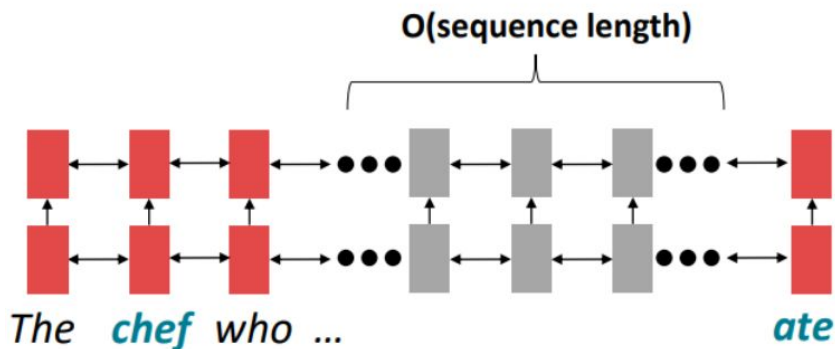
Attention scores

- dot-product attention: $\hat{w}_j = \langle s_j, t \rangle$ (Is it better to normalize this or rescale it by the dimension factor?)
- multiplicative attention: $\hat{w}_j = \langle s_j, Wt \rangle$
- "additive attention": $\hat{w}_j = v^T \sigma(W_1 s_j + W_2 t)$

The actual weights are attention scores passed through **softmax**

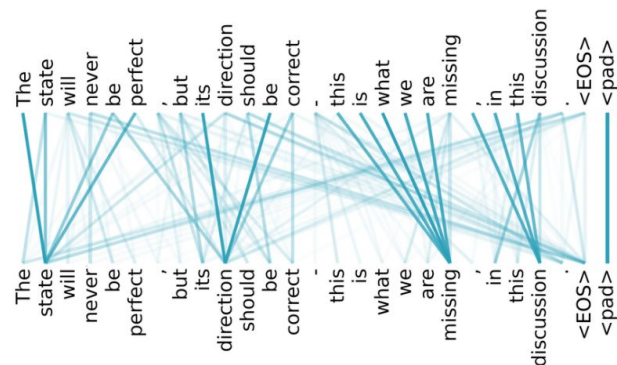
$$w_j = \frac{\exp(\hat{w}_j)}{\sum_k \exp(\hat{w}_k)}$$

Self-attention



RNN

- Long interaction distance
- Resistant to parallelization



Self-attention

- $O(1)$ interaction distance
- Highly parallelizable

Each token gets a selective summary of information from all others

Self-attention

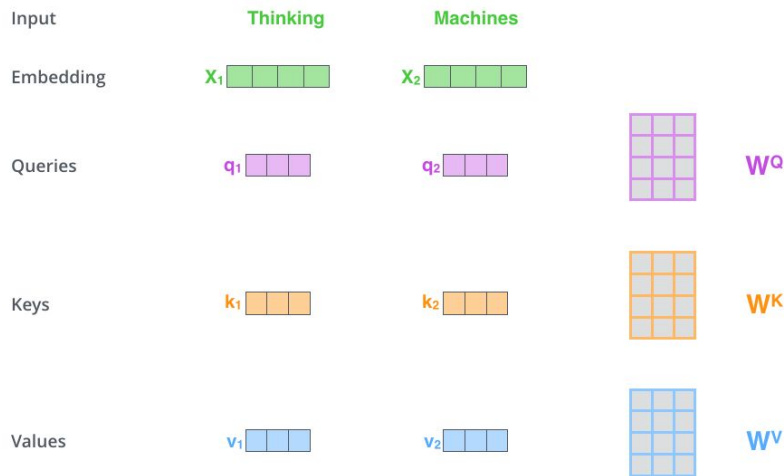


Image credit: <https://jalammargithub.io/illustrated-transformer/>

- Each word now encoded as (query, key, value) triple
- For an input x_i , we have:

$$q_i = (W^Q)^T x_i, \quad k_i = (W^K)^T x_i, \quad v_i = (W^V)^T x_i$$

- Calculate attention scores between query and all keys: $e_{ij} = \langle q_i, k_j \rangle$
- softmax normalization $w_{ij} = \exp(e_{ij}) / \sum_k \exp(e_{ik})$
- output the weighted sum of values $\sum_j w_{ij} v_j$

In matrix notation

- Compute queries, keys, and values

$$Q = XW^Q, \quad K = XW^K, \quad V = XW^V$$

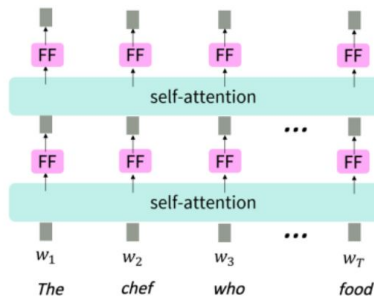
- Calculate attention scores between query and all keys: $E = QK^T$
- softmax normalization $A = \text{softmax}(E)$
- output the weighted sum of values AV

$$\text{output} = \text{softmax}(QK^T)V$$

Question: why we need both query and key?

Equation for Feed Forward Layer

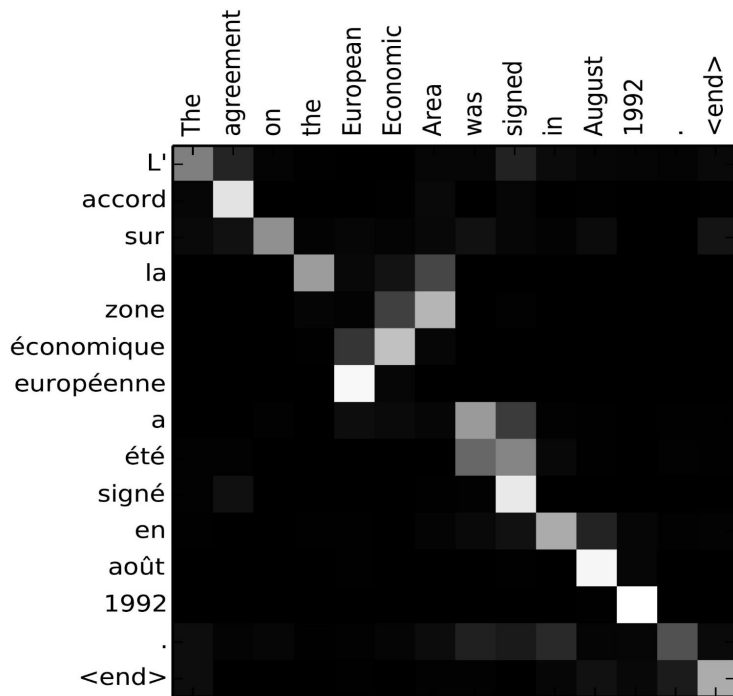
$$m_i = \text{MLP}(\text{output}_i) \\ = W_2 * \text{ReLU}(W_1 \times \text{output}_i + b_1) + b_2$$



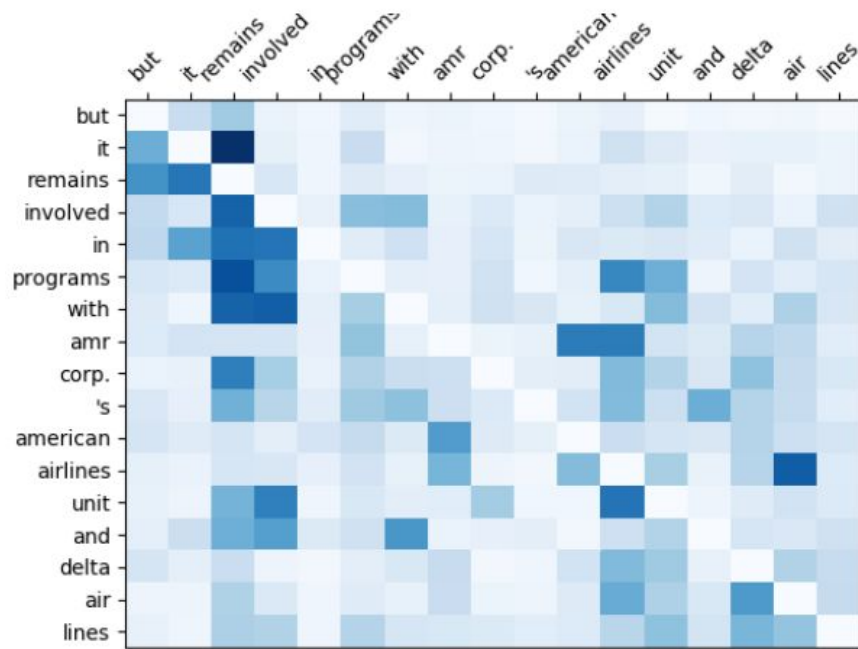
Adding in nonlinearity!

First step toward Transformers!

Attention matrices—visualizing correlations



General attention



Self-attention A

Transformers

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

Illia Polosukhin* ‡
illia.polosukhin@gmail.com

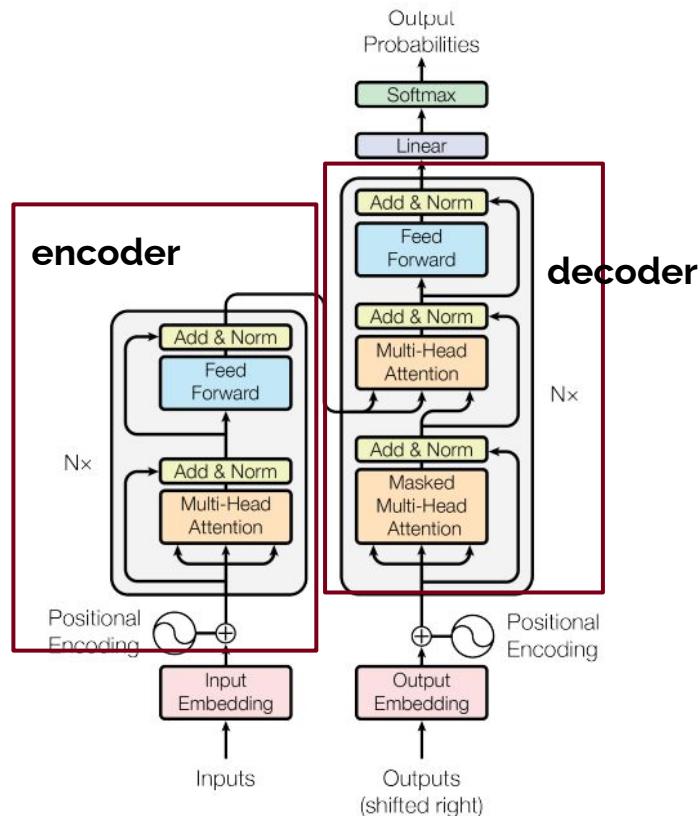


Figure 1: The Transformer - model architecture.

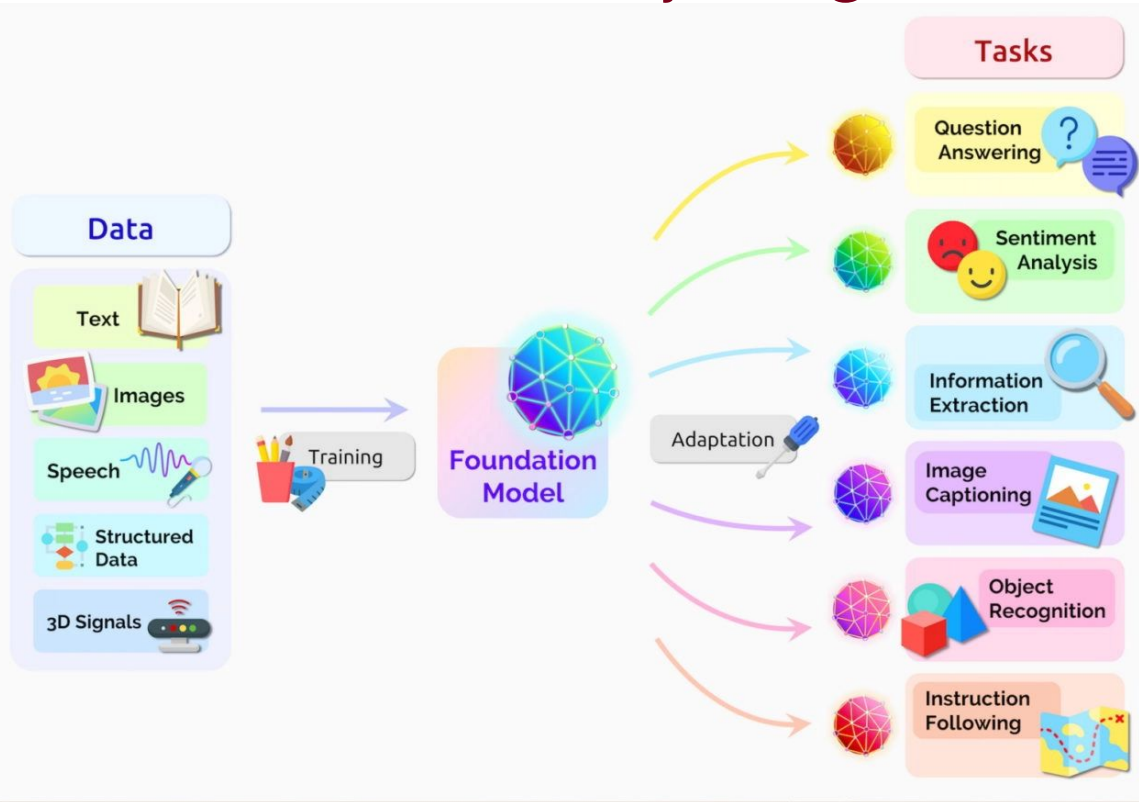
NIPS 2017; <https://arxiv.org/abs/1706.03762>

Transformers reign in NLP!



Image credit: <https://medium.com/mlearning-ai/evolution-of-transformers-part-1-faac3f19d780>

Transformers for everything!

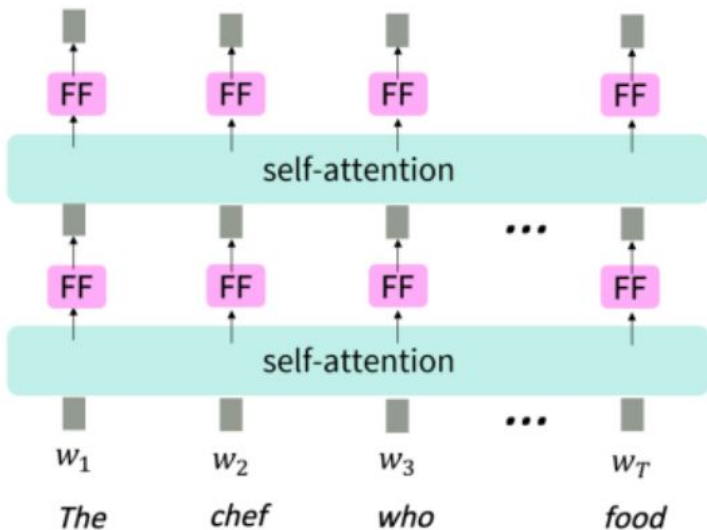


- Transformers have been modified to deal with **almost all** kinds of structured and **unstructured** data
- Enable multimodal data integration and interaction

Starting from self-attention

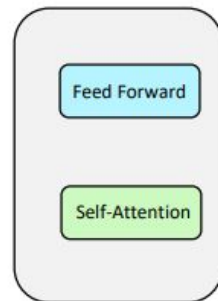
Equation for Feed Forward Layer

$$m_i = MLP(\text{output}_i) \\ = W_2 * \text{ReLU}(W_1 \times \text{output}_i + b_1) + b_2$$



(Credit: Stanford CS231N)

Encoder



Input
Embedding

Inputs



Output
Embedding

Outputs
(shifted right)

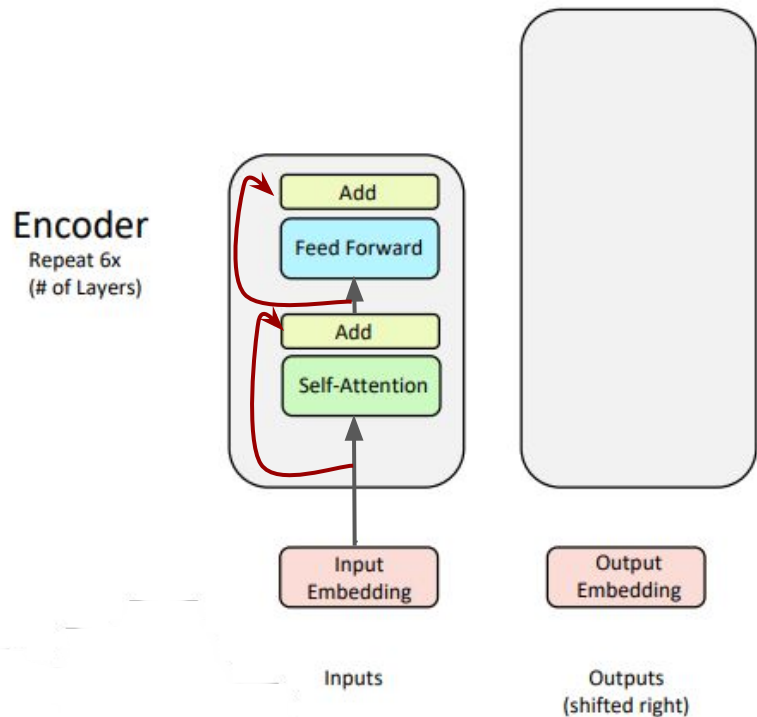
Decoder

(Credit: Stanford CS231N)

Three tricks to build in depth:

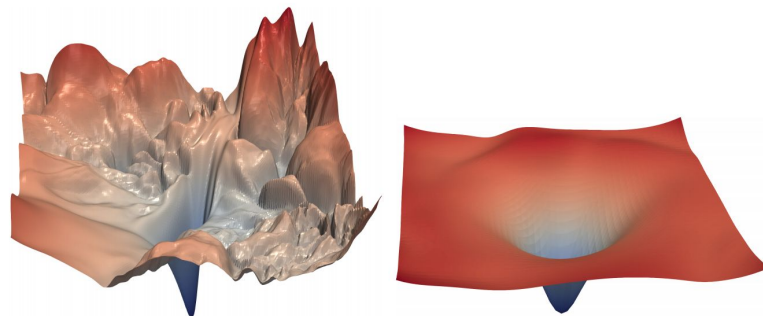
- Residual connection
- Layer normalization
- Scaled inner product attention

Trick 1: Residual connection



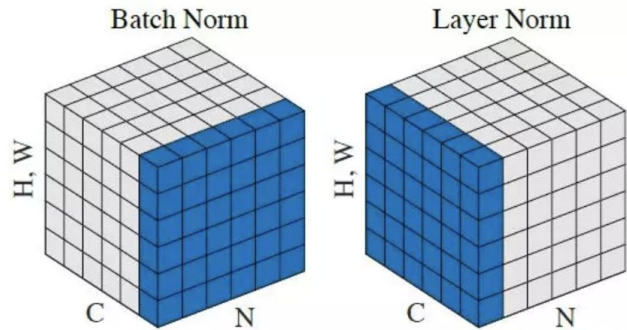
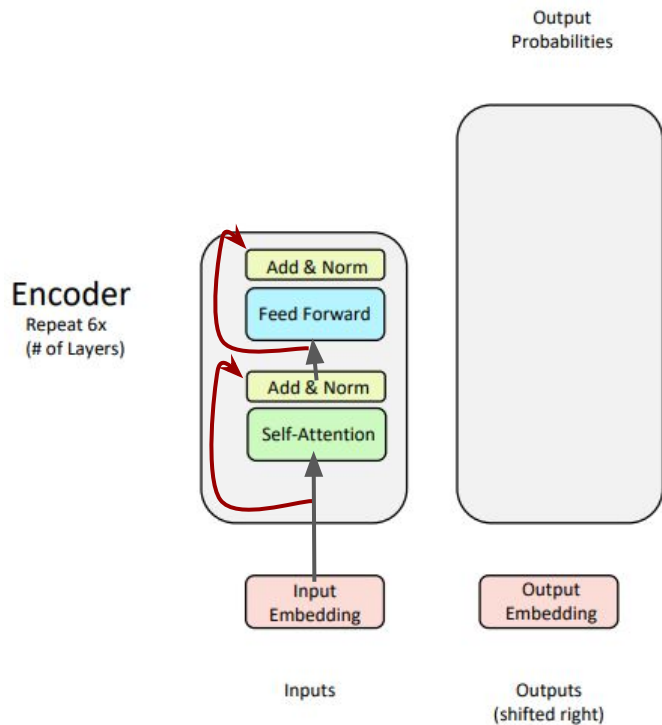
$$\mathbf{x}_k = F(\mathbf{x}_{k-1}) + \mathbf{x}_{k-1}$$

- Mitigating vanishing gradient
- Smoothing out landscape



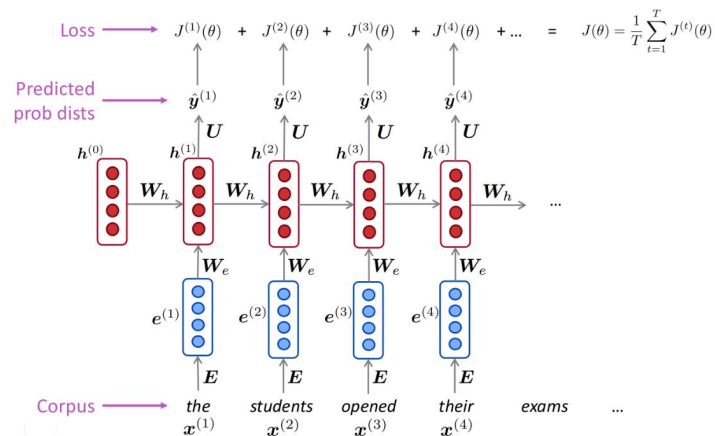
<https://arxiv.org/abs/1712.09913>

Trick 2: Layer normalization

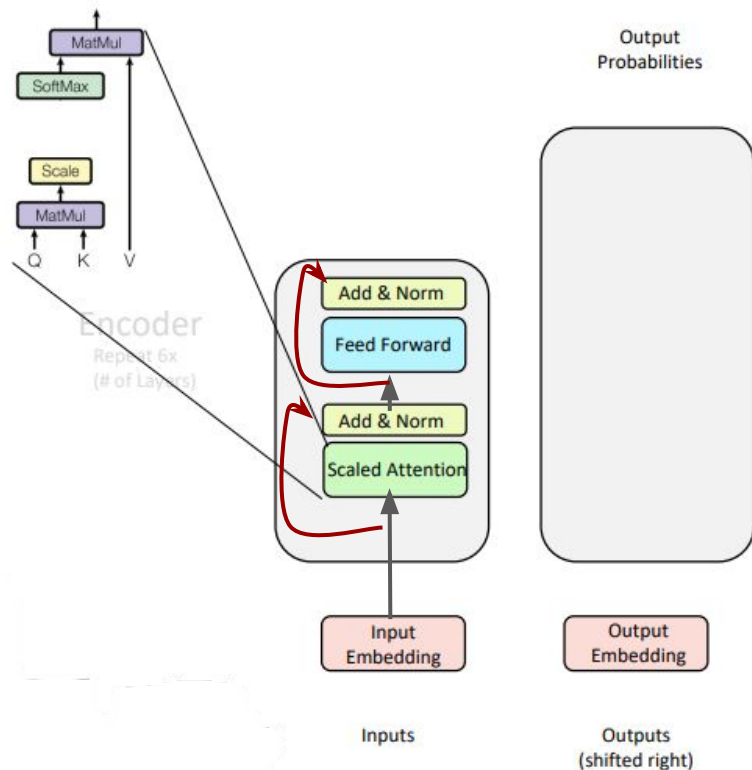


$$x^{\ell'} = \frac{x^{\ell} - \mu^{\ell}}{\sigma^{\ell} + \epsilon}$$

Why not batchnorm?



Trick 3: Scaled inner product attention



$$\text{output} = \text{softmax}(\mathbf{Q}\mathbf{K}^T)\mathbf{V}$$

- After LayerNorm, entries of \mathbf{Q} and \mathbf{K} behaves like IID zero-mean, unit variance
- $\mathbb{E}\langle \mathbf{q}^i, \mathbf{k}^j \rangle = 0$ but

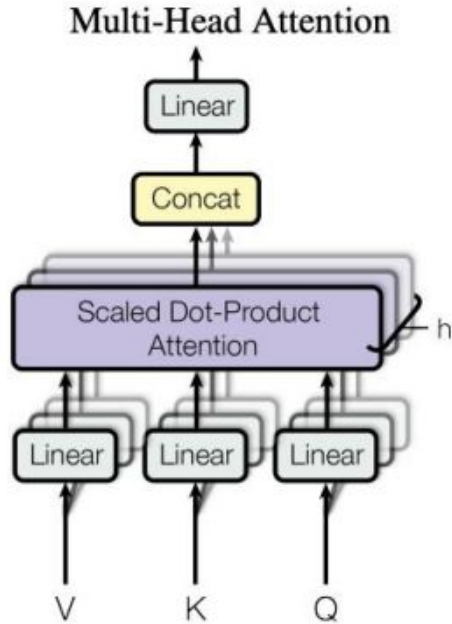
$$\text{Var}\langle \mathbf{q}^i, \mathbf{k}^j \rangle = d_k$$

This can blow up exp computation in the softmax normalization for large d_k !

Solution: normalize by standard deviation

$$\text{output} = \text{softmax}(\mathbf{Q}\mathbf{K}^T / \sqrt{d_k})\mathbf{V}$$

Multi-head attention



[Vaswani et al. 2017]

Multiple, independent self-attention blocks in parallel

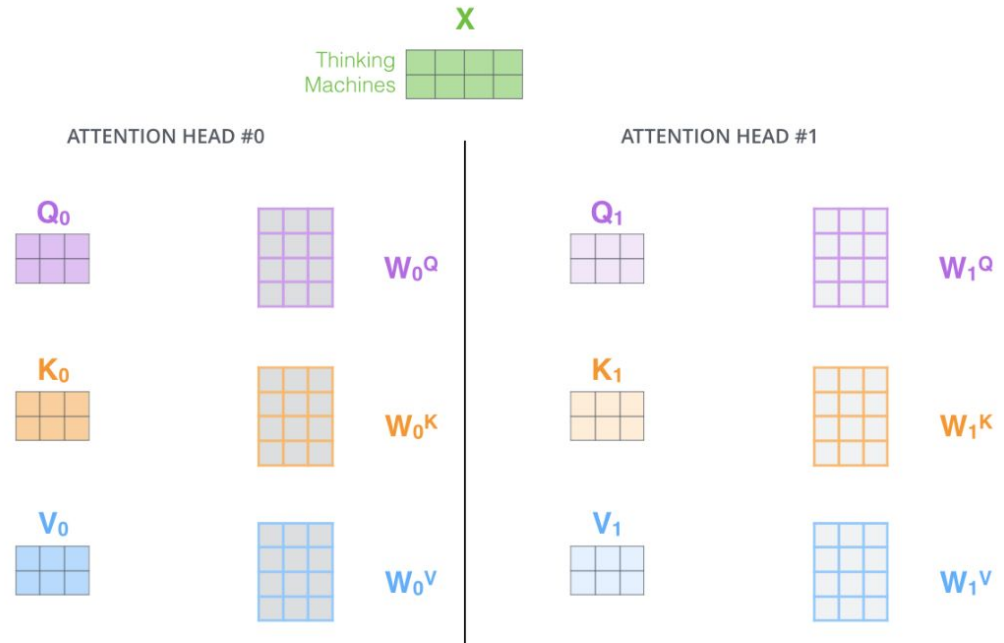
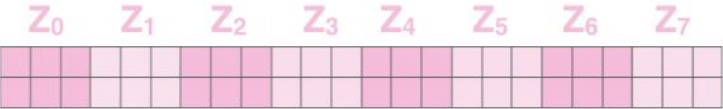


Image credit: <https://ialammargithub.io/illustrated-transformer/>

Intuition: allow the flexibility of capturing different kinds of “relevance”/correlations

Multi-head attention

1) Concatenate all the attention heads



Concatenate

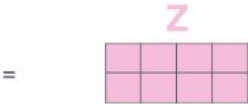


2) Multiply with a weight matrix W^O that was trained jointly with the model

X

Multiply

3) The result would be the Z matrix that captures information from all the attention heads. We can send this forward to the FFNN



Output



Image credit: <https://jalammar.github.io/illustrated-transformer/>

Multi-head attention

1) This is our input sentence*

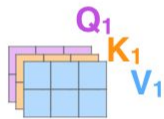
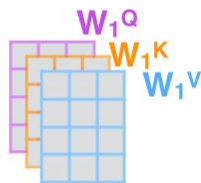
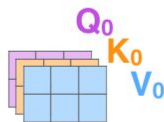
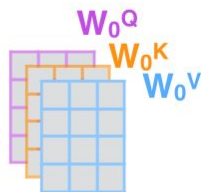
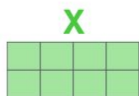
2) We embed each word*

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

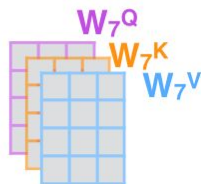
Thinking Machines



...

...

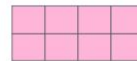
...



W^O

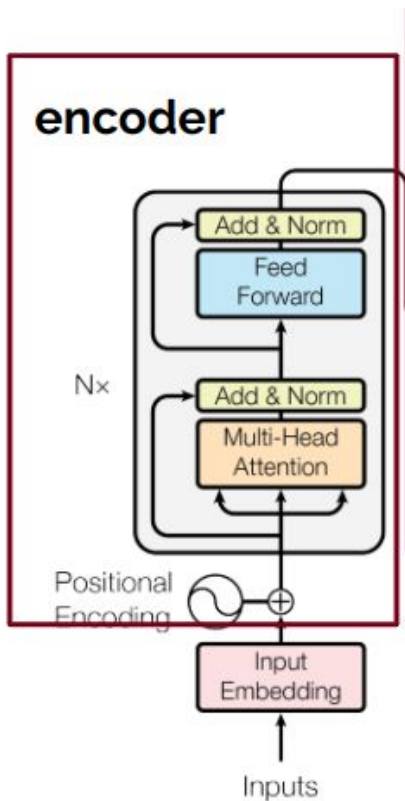


Z



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

Positional encoding



Does the input order matter or not?

$$Q = XW^Q, \quad K = XW^K, \quad V = XW^V$$

$$\text{output} = \text{softmax}(QK^T / \sqrt{d_k})V$$

Positional encoding to break the **order invariance**

- Idea: a positional vector to (hopefully) encode the position information

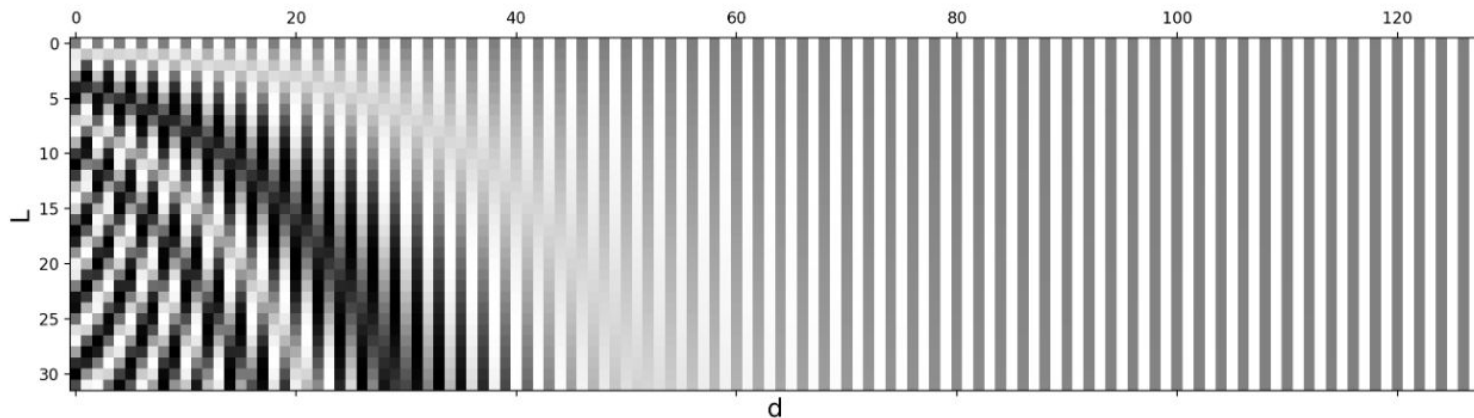
E.g., $X_p = X + P$, or $X_p = [X, P]$

- P can be pre-defined, or made learnable

Sinusoidal positional encoding

$$\text{PE}(i, \delta) = \begin{cases} \sin\left(\frac{i}{10000^{2\delta'/d}}\right) & \text{if } \delta = 2\delta' \\ \cos\left(\frac{i}{10000^{2\delta'/d}}\right) & \text{if } \delta = 2\delta' + 1 \end{cases}$$

L: sequence length
d: embedding dimension



Decoder

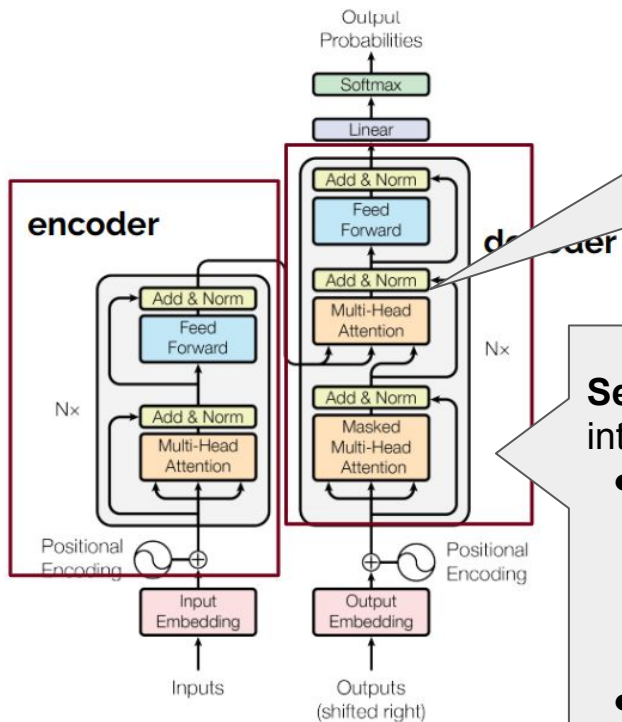
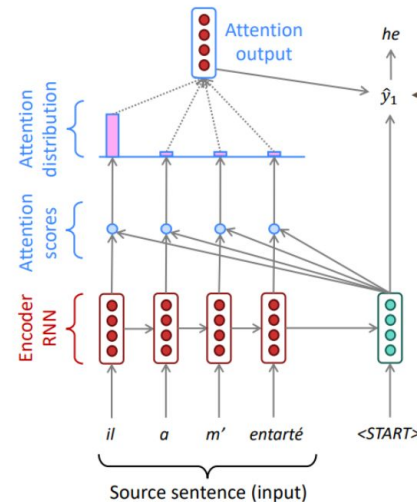


Figure 1: The Transformer - model architecture.

Cross-attention (to model the interaction between the encoder key-values and the current decoder query)

Self-attention (to model the interaction within itself)

- Respect the sequential nature (e.g., language modeling, assuming access to the future is cheating!)
- Masked out future tokens



(Credit: Stanford CS231N)

We can look at these (not greyed out) words

For encoding these words

	[START]	The	chef	who
[START]	-∞	-∞	-∞	-∞
The		-∞	-∞	-∞
chef			-∞	-∞
who				-∞

Attention matrix

Computation

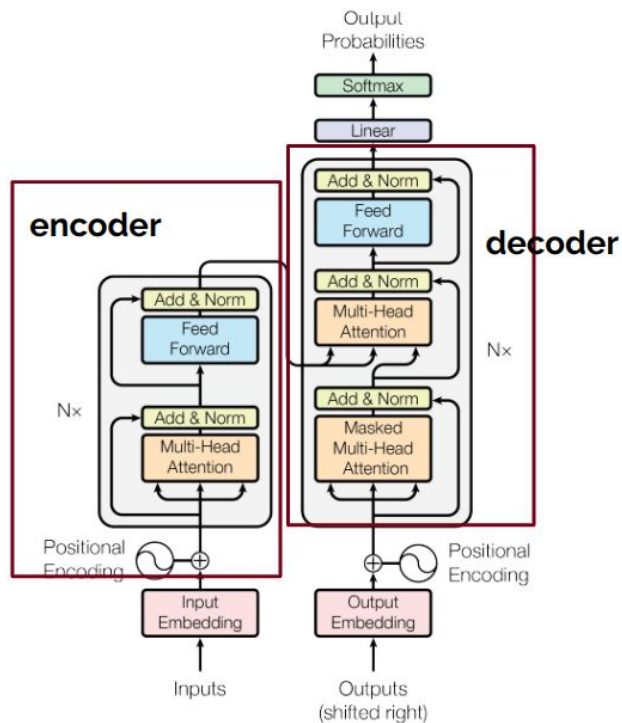


Figure 1: The Transformer - model architecture.

What's the total computation?

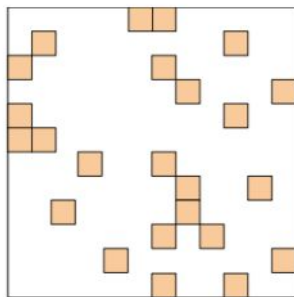
$$Q = XW^Q, \quad K = XW^K, \quad V = XW^V$$

$$\text{output} = \text{softmax}(QK^T / \sqrt{d_k})V$$

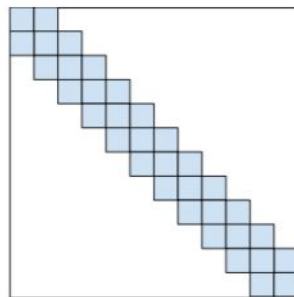
$$O(T^2 d)$$

Quadratic computation vs. linear computation in RNNs

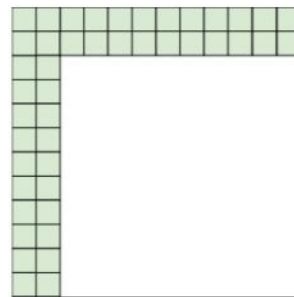
Computation



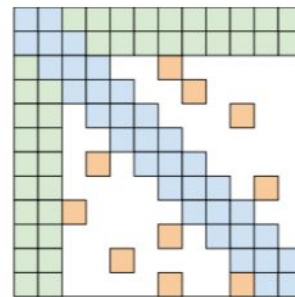
(a) Random attention



(b) Window attention



(c) Global Attention



(d) BIGBIRD

Idea; building in sparsity <https://arxiv.org/abs/2007.14062>

Do Transformer Modifications Transfer Across Implementations and Applications?

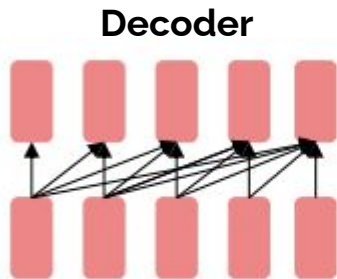
Sharan Narang*	Hyung Won Chung	Yi Tay	William Fedus
Thibault Fevry†	Michael Matena†	Karishma Malkan†	Noah Fiedel
Noam Shazeer	Zhenzhong Lan†	Yanqi Zhou	Wei Li
Nan Ding	Jake Marcus	Adam Roberts	Colin Raffel†

But not much consistent improvement so far

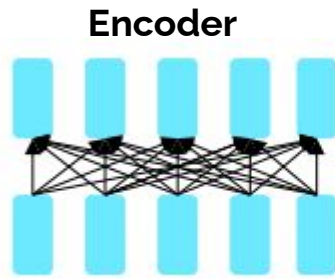
<https://arxiv.org/abs/2102.11972>

Pretraining

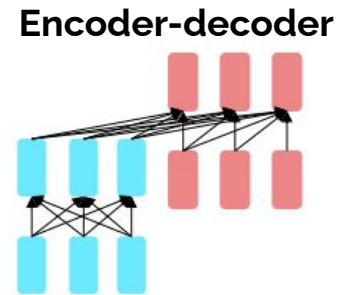
Pretraining + finetuning pipeline is standard in modern NLP/CV and many applied areas



Good for content generation
(e.g., GPT-2, GPT-3)

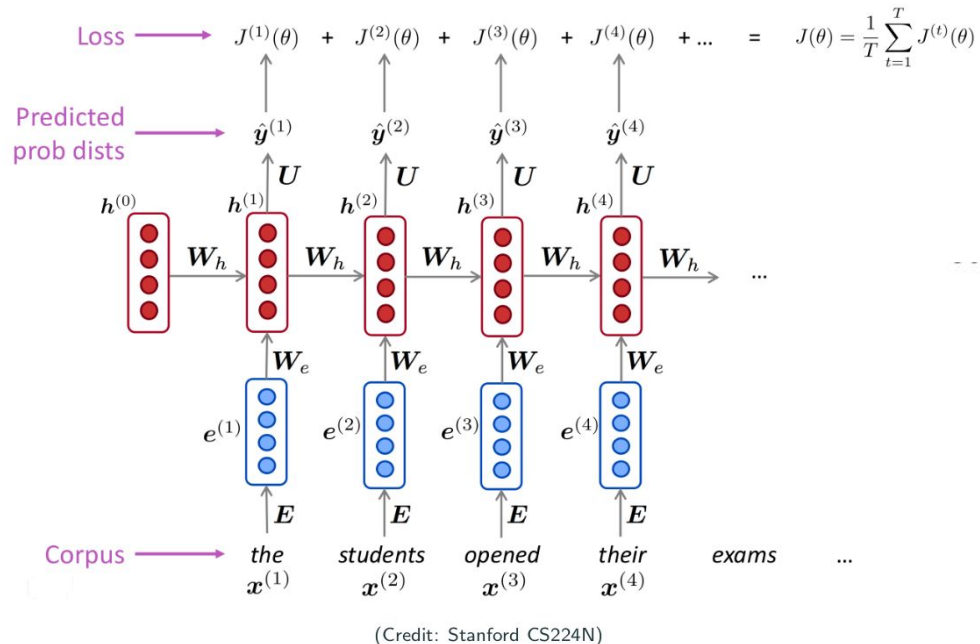


Good for feature extraction
(e.g., X-**BERT**)

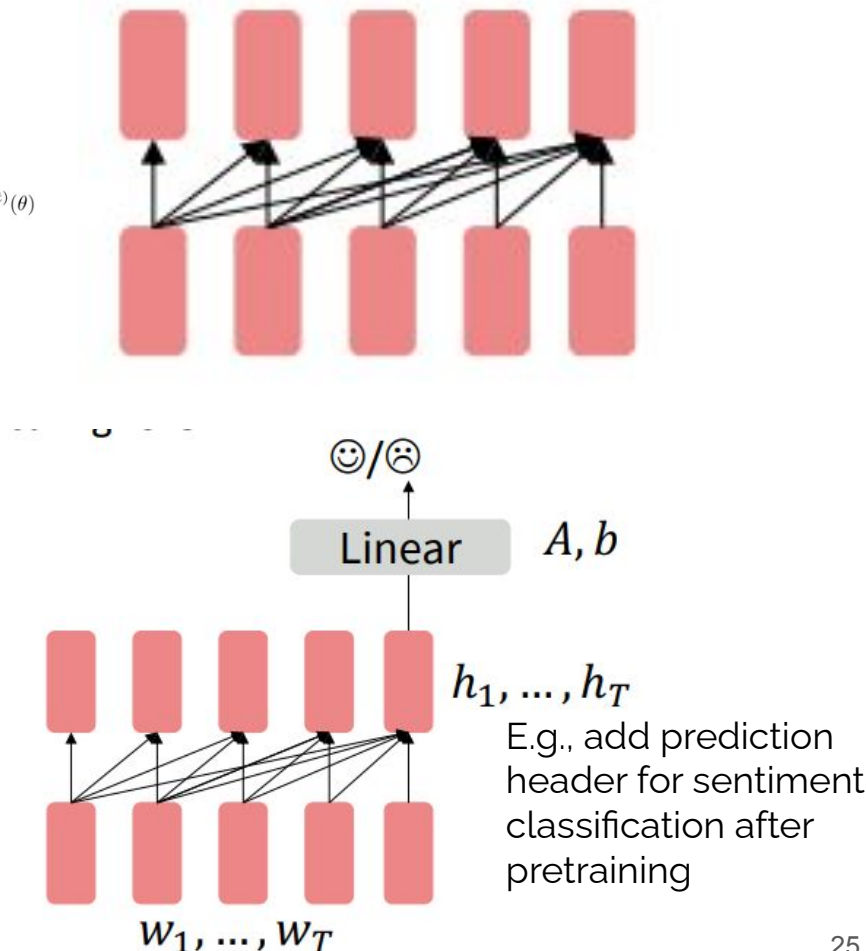


Good for everything?
(e.g., Transformer, T5)

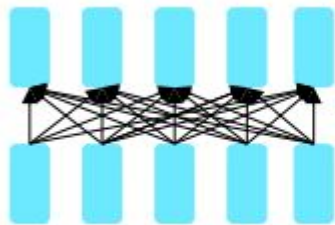
Decoder pretraining



Can be pretrained on language modeling, i.e., model $\mathbb{P}[x^{(1)}, \dots, x^{(T)}]$ or $\mathbb{P}[x^{(t+1)} | x^{(t)}, \dots, x^{(1)}]$



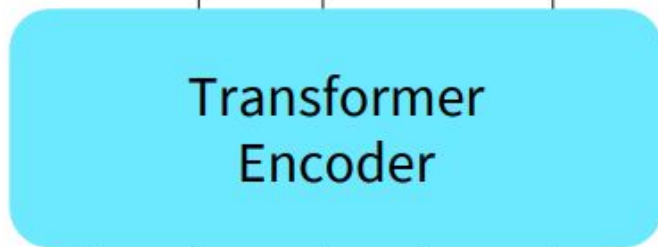
Encoder pretraining



Both past and future information available due to the bi-direction modeling; not ideal for language modeling

[Predict these!]

went to store



I pizza to the [M]

[Replaced] [Not replaced] [Masked]

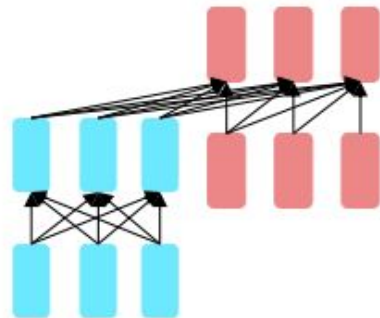
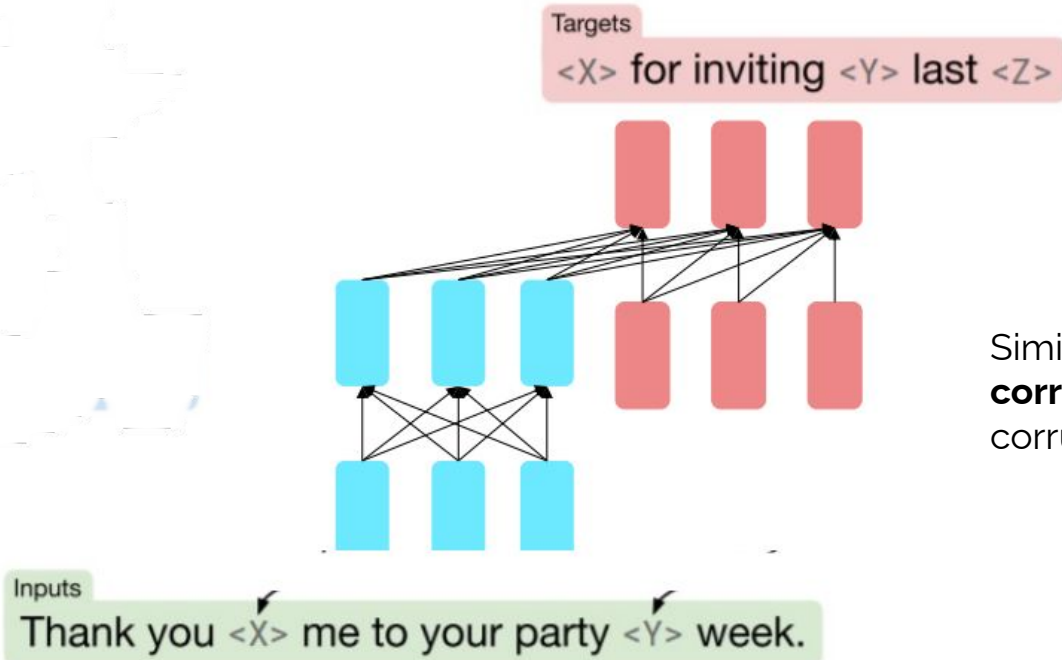
Idea: create corruption and predict the right things

- Masked out words (i.e., missing words)
- Noisy words (randomly replaced)

Self-supervised learning! (pretext tasks)

<https://arxiv.org/abs/1810.04805>

Encoder-decoder pretraining



Similar to pretraining encoder,
corruption removal! (called span corruption)

<https://arxiv.org/abs/1910.10683>