

Recurrent Neural Networks and Trasformers in Natural Language Processing

Outline

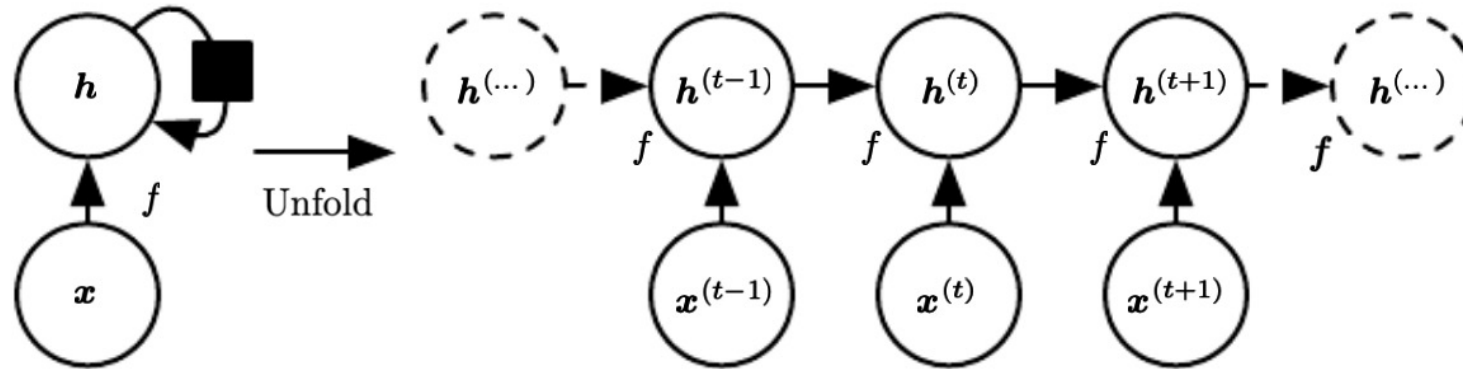
- Recurrent Neural Networks
- Long Short-Term Memory
- Transformers
- Natural Language Processing
- BERT

Recurrent Neural Networks (RNN)

- Neural network that is specialized for processing a sequence of values $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(\tau)}$
- Each member of the output is a function of the previous members of the output. Each member of the output is produced using the same update rule applied to the previous outputs.

$$\mathbf{s}^{(t)} = f(\mathbf{s}^{(t-1)}; \boldsymbol{\theta})$$

Recurrent Neural Networks (RNN)

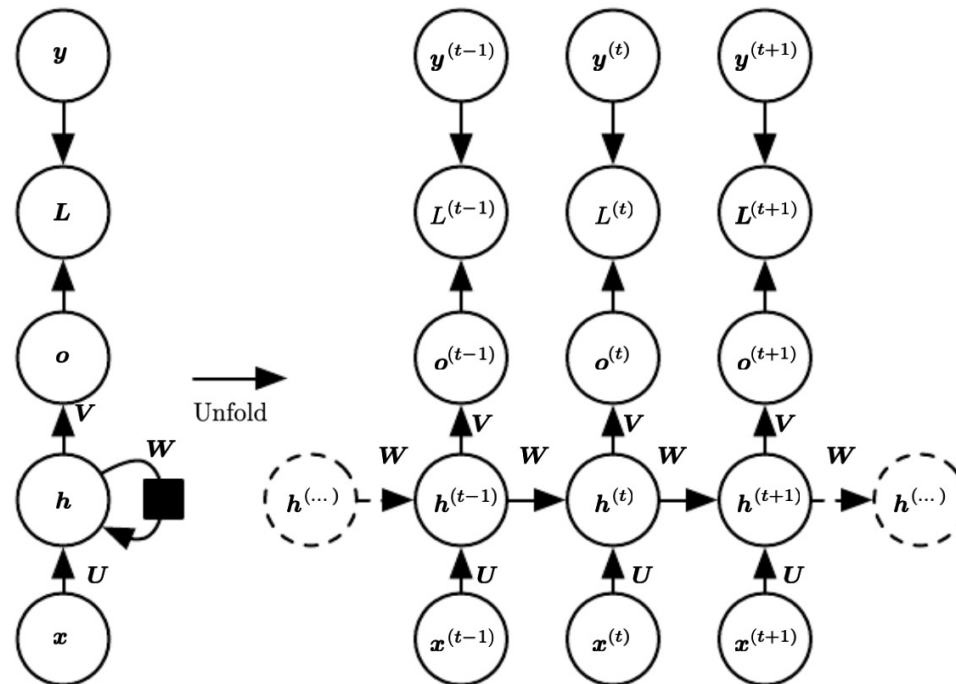


$$\begin{aligned}h^{(t)} &= f(h^{(t-1)}, x^{(t)}; \theta) \\h^{(t-1)} &= f(h^{(t-2)}, x^{(t-1)}; \theta) \\&\vdots \\h^{(1)} &= f(h^{(0)}, x^{(1)}; \theta)\end{aligned}$$

Recurrent Neural Networks (RNN)

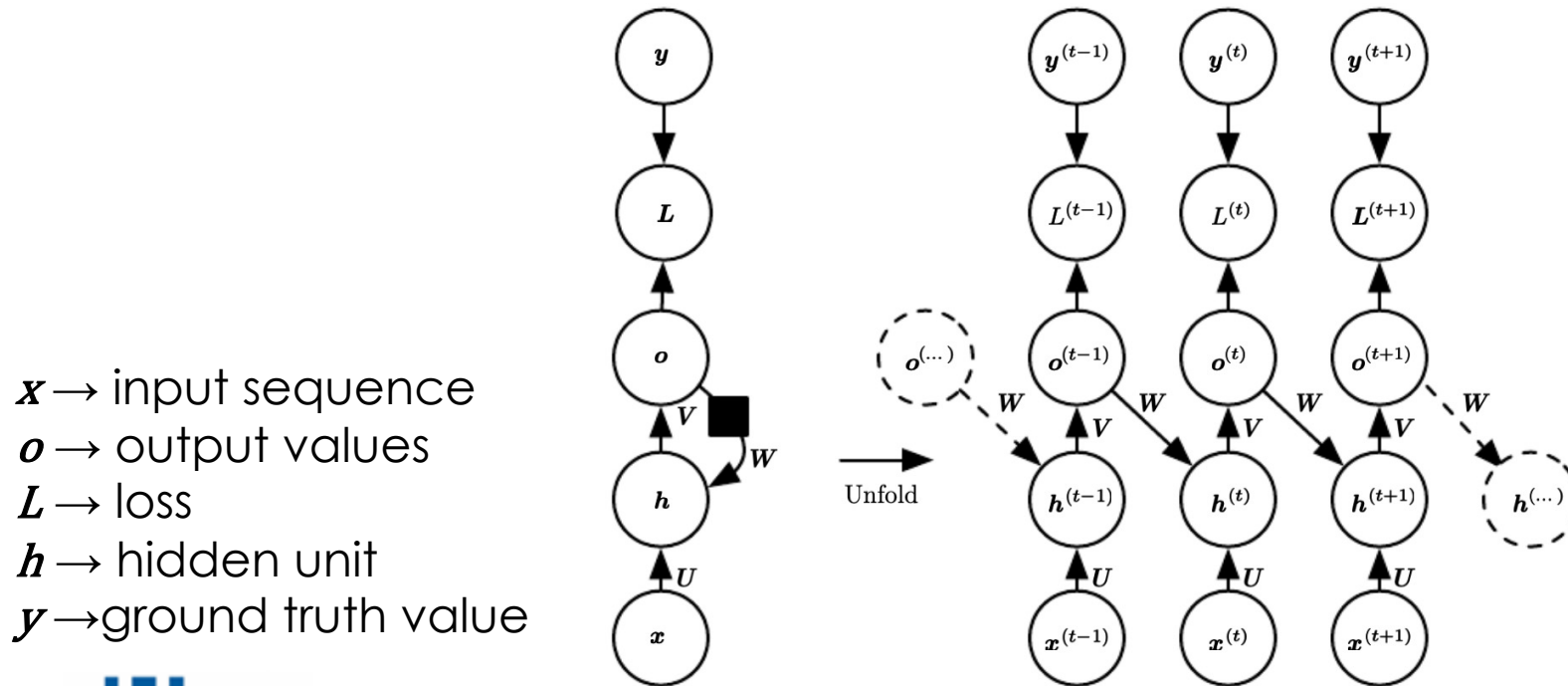
Recurrent networks that produce an output at each time step and have recurrent connections between hidden units.

$x \rightarrow$ input sequence
 $o \rightarrow$ output values
 $L \rightarrow$ loss
 $h \rightarrow$ hidden unit
 $y \rightarrow$ ground truth value



Recurrent Neural Networks (RNN)

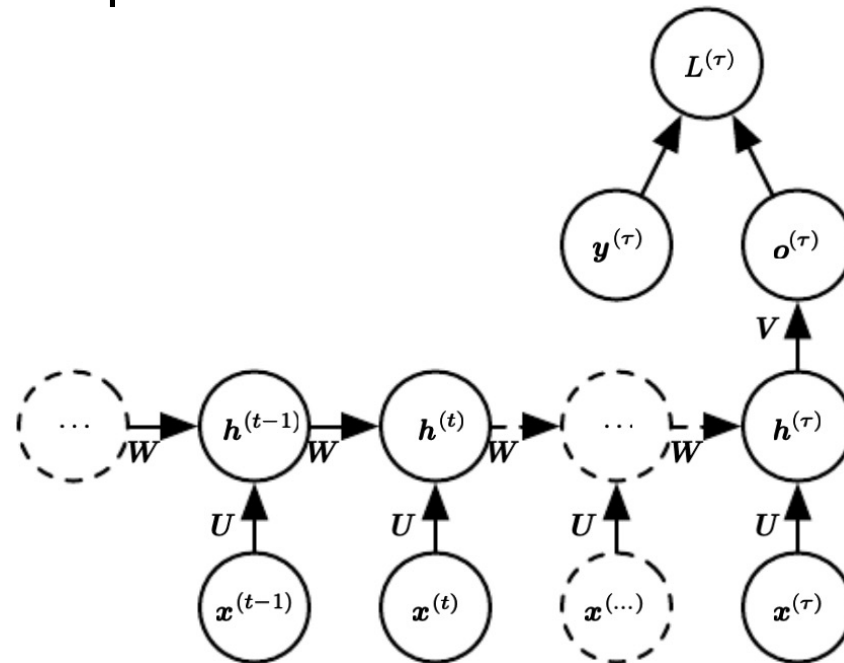
Recurrent networks that produce an output at each time step and have recurrent connections only from the output at one time step to the hidden units at the next time step.



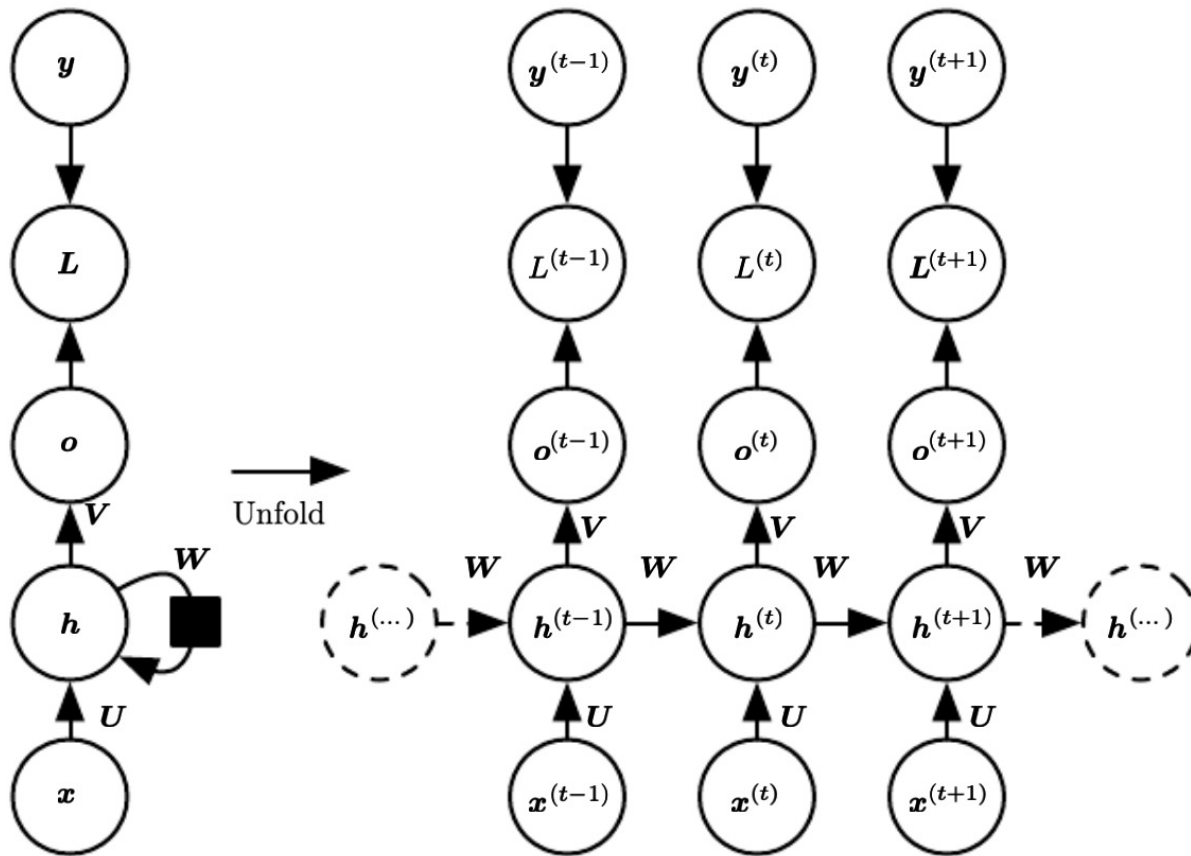
Recurrent Neural Networks (RNN)

Recurrent networks with recurrent connections between hidden units, that read an entire sequence and then produce a single output.

$x \rightarrow$ input sequence
 $o \rightarrow$ output values
 $L \rightarrow$ loss
 $h \rightarrow$ hidden unit
 $y \rightarrow$ ground truth value



Recurrent Neural Networks (RNN)



$$\mathbf{a}^{(t)} = \mathbf{b} + \mathbf{W}\mathbf{h}^{(t-1)} + \mathbf{U}\mathbf{x}^{(t)}$$

$$\mathbf{h}^{(t)} = \tanh(\mathbf{a}^{(t)})$$

$$\mathbf{o}^{(t)} = \mathbf{c} + \mathbf{V}\mathbf{h}^{(t)}$$

$$\hat{\mathbf{y}}^{(t)} = \text{softmax}(\mathbf{o}^{(t)})$$

$$\begin{aligned} L(\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(\tau)}\}, \{\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(\tau)}\}) &= \\ &= \sum_t L^{(t)} = \\ &= - \sum_t \log p_{\text{model}}(\mathbf{y}^{(t)} | \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}\}) \end{aligned}$$

Recurrent Neural Networks (RNN)

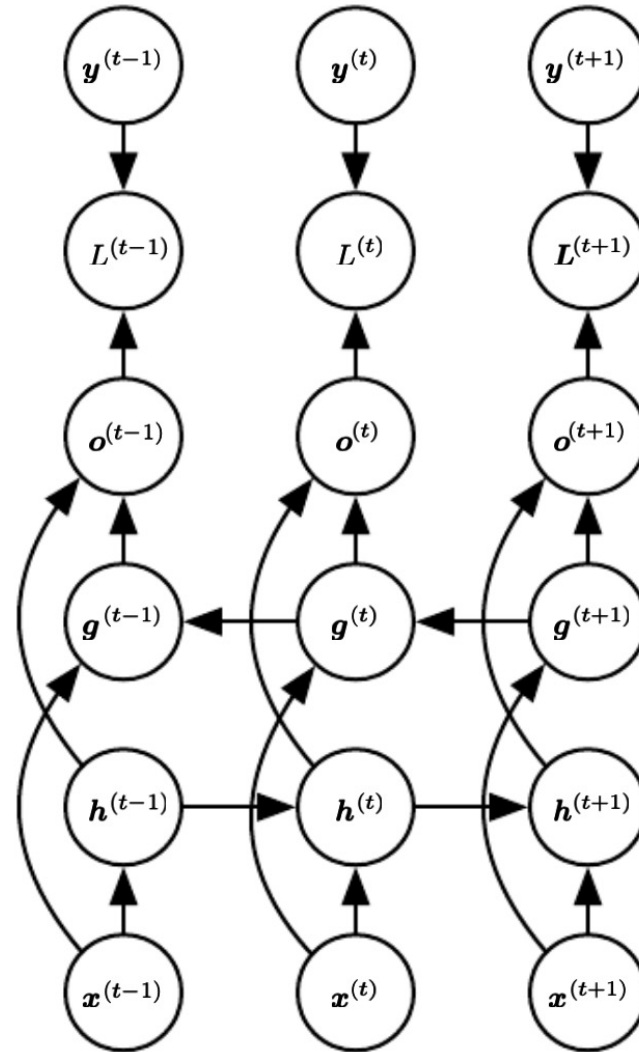
Training:

- Back-propagation through time (BPTT)
- Teacher forcing

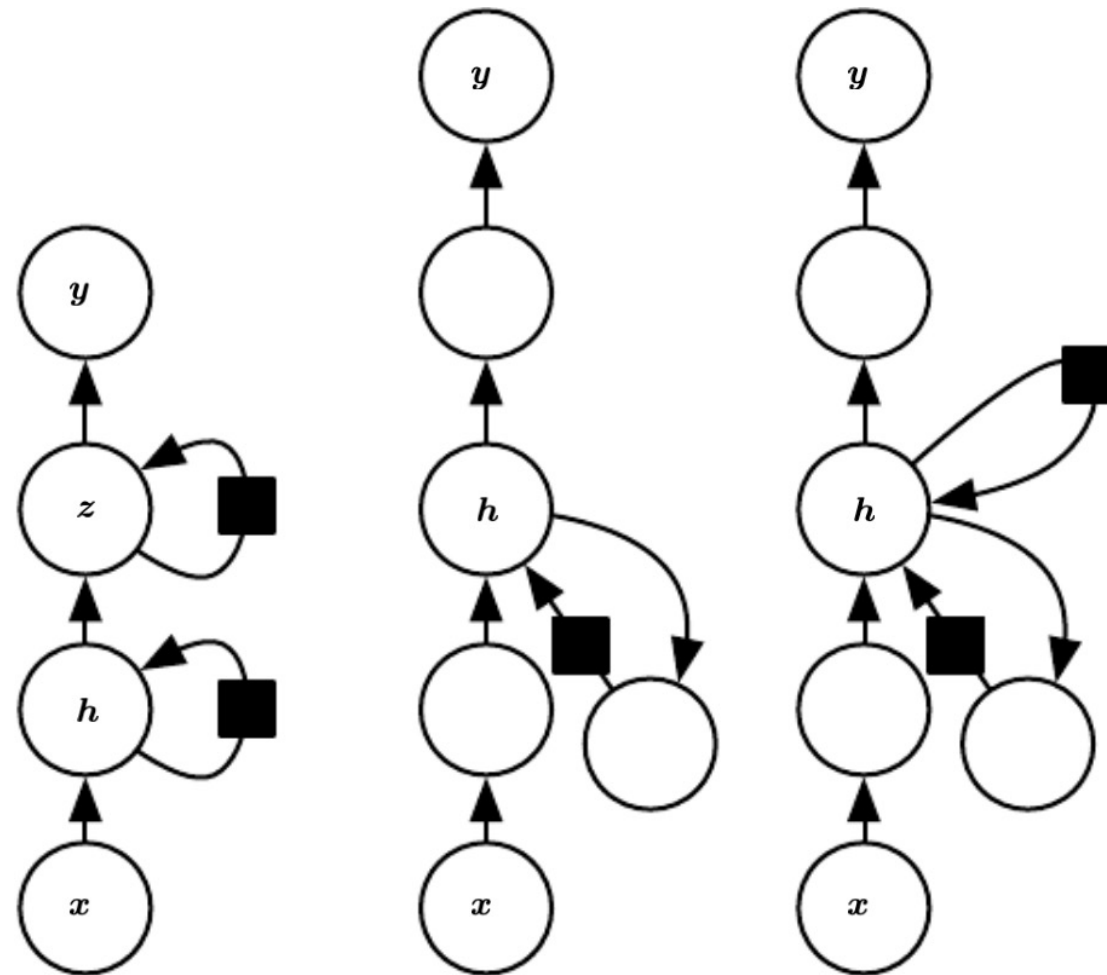
Difficult in training

- Attention window
- Reduced number of parameters

BiDirectional RNN

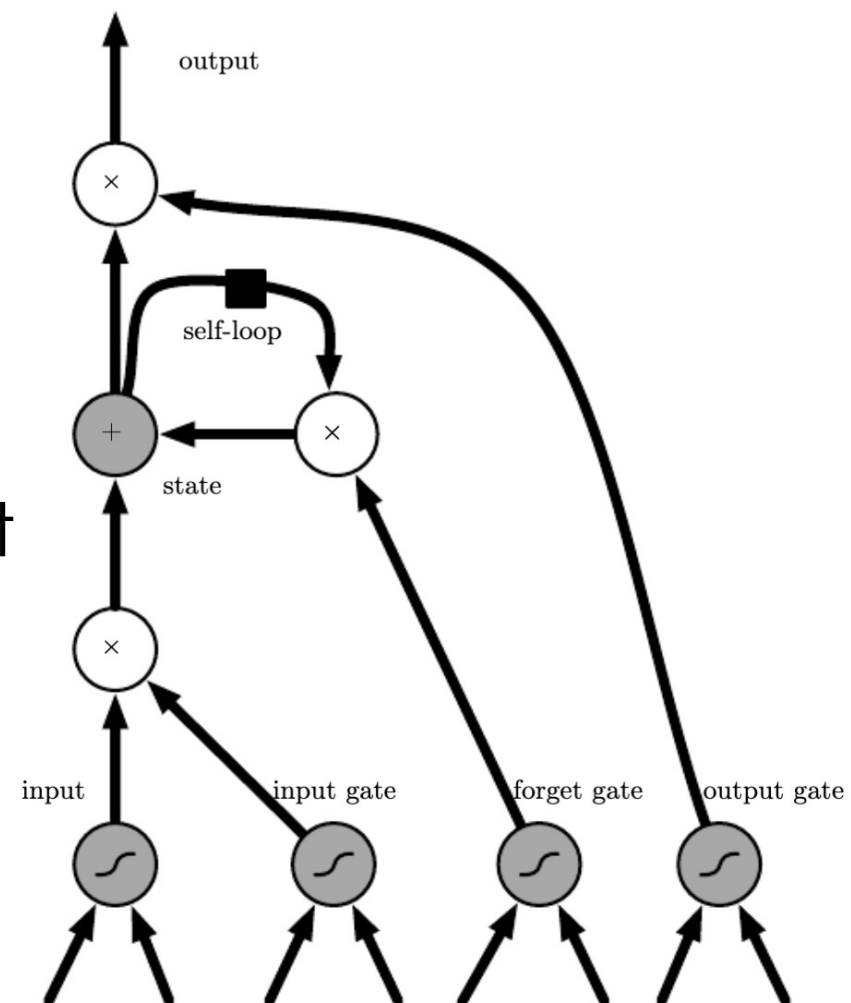


Deep Recurrent Networks



Long Short-Term Memory (LSTM)

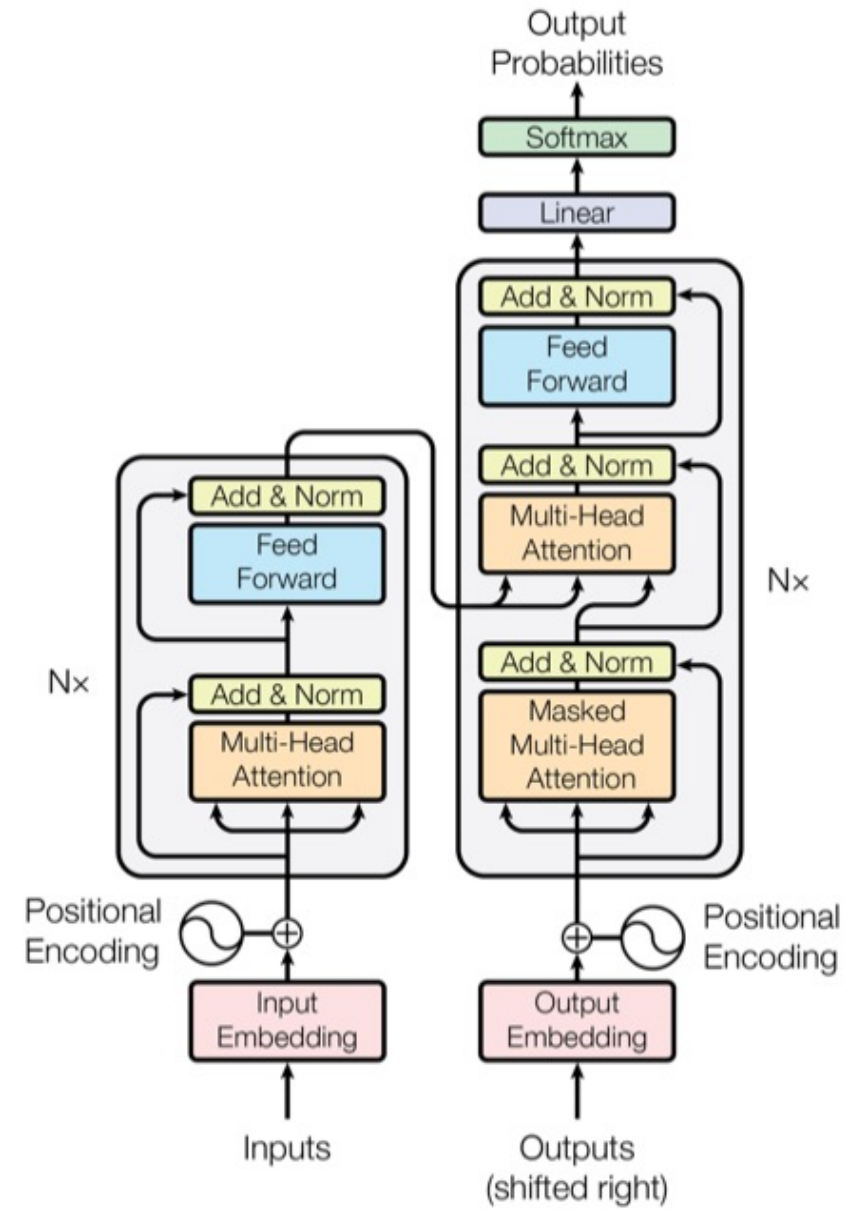
LSTM learn long-term dependencies more easily than the simple recurrent architectures.



Transformers

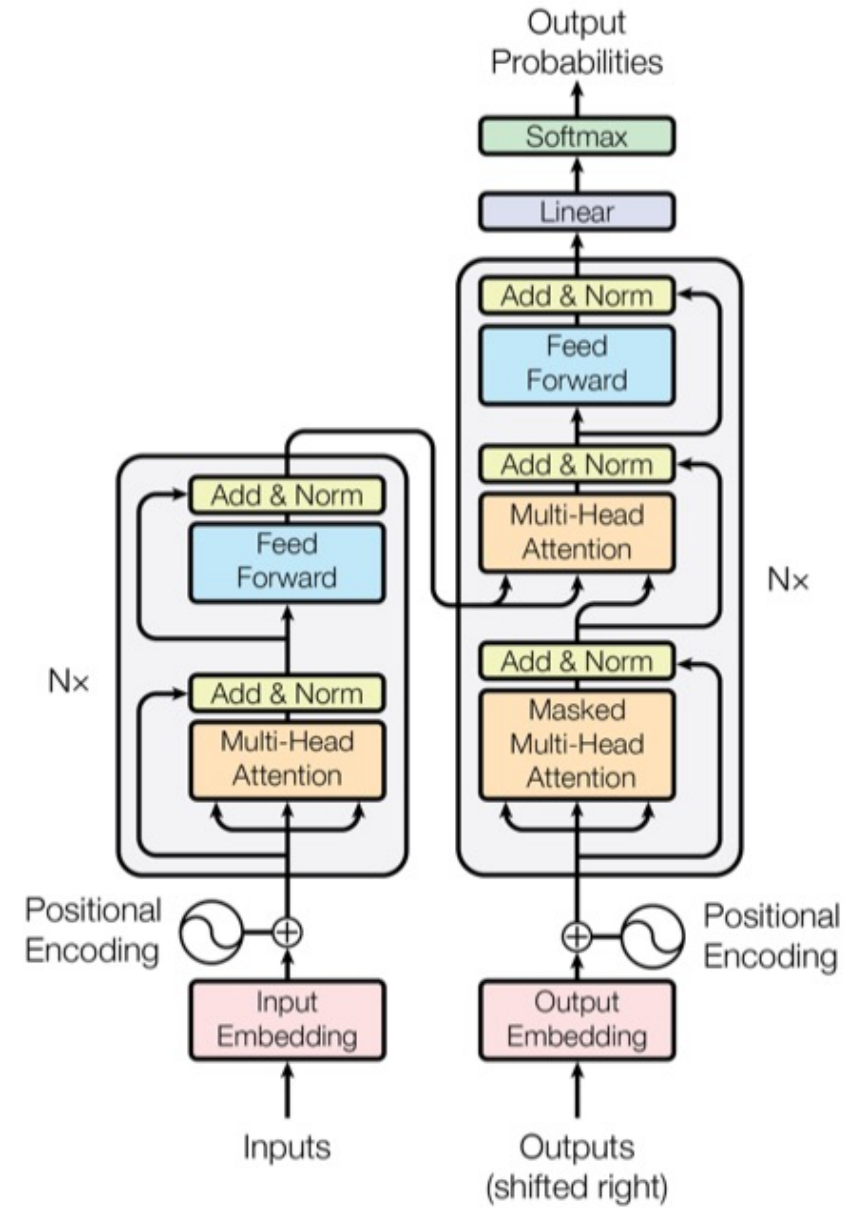
Model architecture delete recurrence and relying entirely on an attention mechanism to draw global dependencies between input and output.

$$\begin{aligned}(\mathbf{z}_1, \dots, \mathbf{z}_n) &= \text{encoder}(\mathbf{x}_1, \dots, \mathbf{x}_n) \\ (\mathbf{y}_1, \dots, \mathbf{y}_m) &= \text{decoder}(\mathbf{z}_1, \dots, \mathbf{z}_n)\end{aligned}$$



Transformers

- Encoder
 - Multi-head self-attention mechanism + residual connection + normalization layer
 - Feed-forward fully connected layer + residual connection + normalization layer
- Decoder
 - Multi-head attention over the output of the encoder stack + residual connection + normalization layer
 - Multi-head self-attention mechanism + residual connection + normalization layer
 - Feed-forward fully connected layer + residual connection + normalization layer



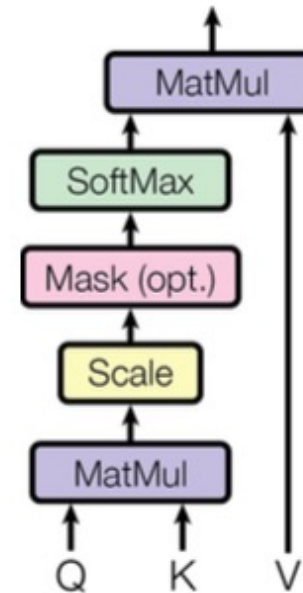
Transformers

From an input embedding, three roles are played:

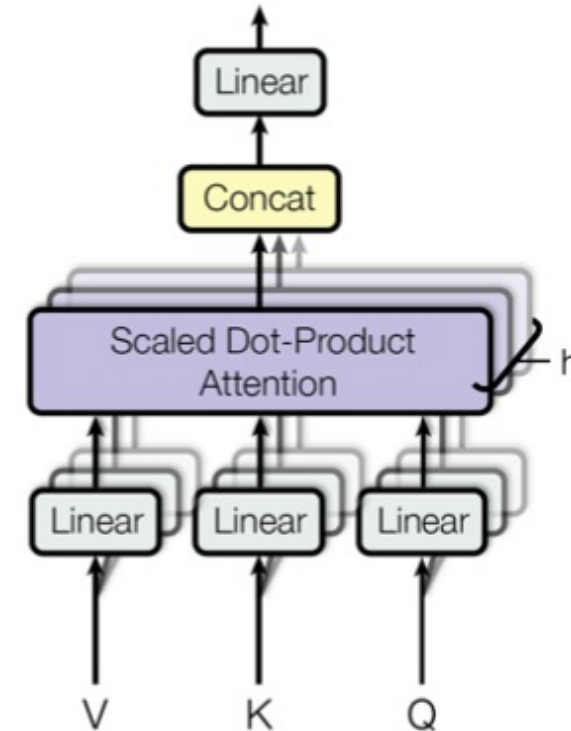
- Query → the current focus of attention
- Key → a preceding input
- Value → the value for the computation

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

Scaled Dot-Product Attention



Multi-Head Attention



Natural Language Processing (NLP)

- NLP is the use of human languages, such as English or French, by a computer.
- Many NLP applications are based on language models that define a probability distribution over sequences of words, characters, or bytes in a natural language.

Natural Language Processing (NLP)

n -grams is a simple language model where the conditional probability of the n -th token is given from the preceding $n-1$ tokens

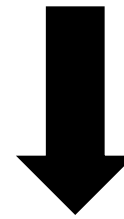
$$P(x_1, \dots, x_\tau) = P(x_1, \dots, x_{n-1}) \prod_{t=n}^{\tau} P(x_t | x_{t-n+1}, \dots, x_{t-1})$$

$$P(x_t | x_{t-n+1}, \dots, x_{t-1}) = \frac{P_n(x_{t-n+1}, \dots, x_t)}{P_{n-1}(x_{t-n+1}, \dots, x_{t-1})}$$

$$P(\text{THE DOG RAN AWAY}) = P_3(\text{THE DOG RAN}) \frac{P_3(\text{DOG RAN AWAY})}{P_2(\text{DOG RAN})}$$

Natural Language Processing (NLP)

A misty ridge uprises from the surge

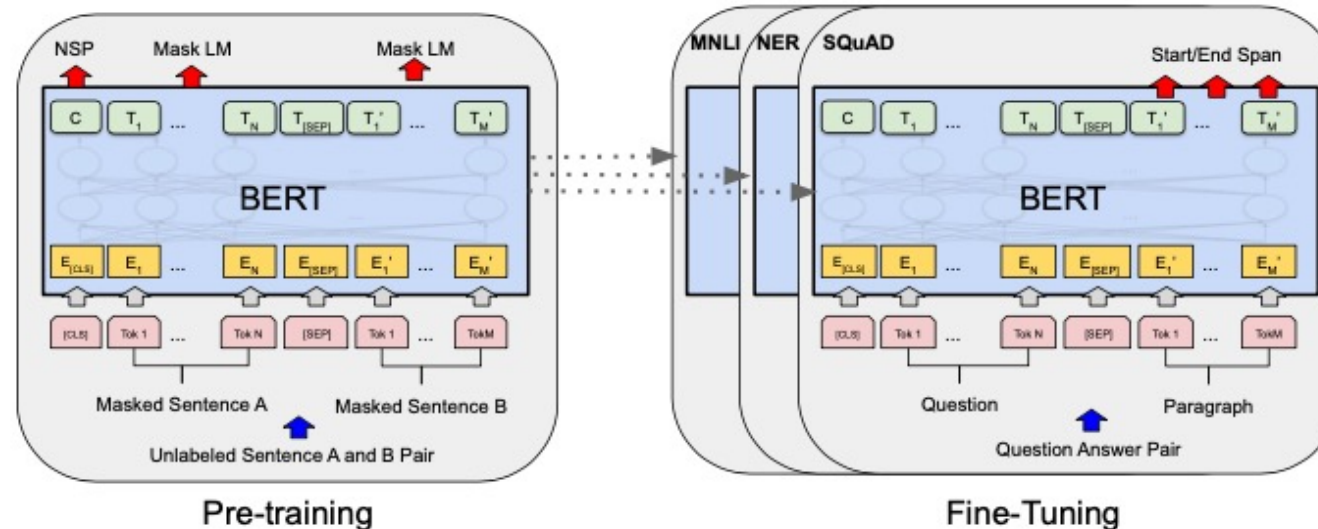


tokenization

['a', 'misty', 'ridge', 'up', '##rise', '##s',
'from', 'the', 'surge']

Bidirectional Encoder Representations from Transformers (BERT)

- Pre-Training
 - Masked Language Model
 - Next Sentence Prediction
- Fine-Tuning



Bidirectional Encoder Representations from Transformers (BERT)

$L \rightarrow$ Number of layers (Transformer blocks)

$H \rightarrow$ Hidden size

$A \rightarrow$ Number of self-attention heads

- BERT_{BASE} ($L=12$, $H=768$, $A=12$, Total Parameters=110M)
- BERT_{LARGE} ($L=24$, $H=1024$, $A=16$, Total Parameters=340M)

