

Sequence Modelling

Mingfei Sun

Foundations of Machine Learning
The University of Manchester



The University of Manchester

A Quick Recap

In the previous lecture, we've covered:

- ▶ Multi-Layer Perceptron (MLP)
- ▶ Feedforward neural networks:
Forward equations, Backpropagation equations
- ▶ Training of neural networks:
Parameter initialization, Normalization, SGD optimization, and Regularization.

Outline

Language Models

Recurrent Neural Networks

Transformers

Intended Learning Outcomes

By the end of this lecture, you'll be able to

- ▶ Identify the key components in language models, Recurrent Neural Networks, and Transformers
- ▶ Explain how the Back-Propagation Through Time (BPTT) and Truncated BPTT can be used to optimize Recurrent Neural Networks
- ▶ Explain how the Attention is used to transform sequences in an equivariant or non-equivariant way
- ▶ Apply and implement self-attention to process sequential inputs

Outline

Language Models

Recurrent Neural Networks

Transformers

Language Models

A language model assigns a probability to a sequence of words, such that

$$\sum_{x \in \Sigma^*} p(x) = 1$$

- ▶ Given the observed training text, how probable is this new utterance?
- ▶ Can compare different orderings of words (e.g. Translation):

$$p(\text{he likes apples}) > p(\text{apples likes he})$$

or choice of words (e.g., Speech Recognition)

$$p(\text{he likes apples}) > p(\text{he licks apples})$$

- ▶ Employ the chain rule to decompose the joint probability into a sequence of conditional probabilities:

$$p(x_1, x_2, x_3, \dots, x_N) = p(x_1)p(x_2|x_1)p(x_3|x_1, x_2) \cdot \dots \cdot p(x_N|x_1, x_2, \dots, x_{N-1})$$

- ▶ Can model complex joint distributions by learning conditional distributions over the next word (x_n) given the history of words observed (x_1, x_2, \dots, x_{n-1})

Language Models

The simple objective of modelling the next word given the observed history contains much of the complexity of natural language understanding.

- ▶ Consider predicting the extension of the utterance:

$$p(\cdot | \text{There she built a })$$

- ▶ With more context we are able to use our knowledge of both language and the world to heavily constrain the distribution over the next word:

$$p(\cdot | \text{Alice went to the beach. There she built a })$$

There is evidence that human language acquisition partly relies on future prediction

- ▶ Language modelling is a time series prediction problem in which we must be careful to train on the past and test on the future

Outline

Language Models

Recurrent Neural Networks

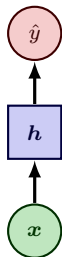
Transformers

Recurrent Neural Networks (RNNs)

Feedforward

$$\mathbf{h} = g(\mathbf{V}\mathbf{x} + \mathbf{c})$$

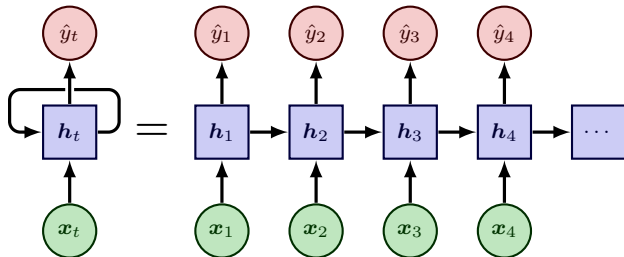
$$\hat{\mathbf{y}} = \mathbf{W}\mathbf{h} + \mathbf{b}$$



Recurrent neural networks

$$\mathbf{h}_n = g(\mathbf{V}[\mathbf{x}_n; \mathbf{h}_{n-1}] + \mathbf{c})$$

$$\hat{\mathbf{y}}_n = \mathbf{W}\mathbf{h}_n + \mathbf{b}$$



The unrolled recurrent network is a directed acyclic graph.

RNNs: Unrolled step-by-step

Consider an example of Next Word Prediction:

$$p(\cdot | \text{There she built a })$$

How to build an RNN model for this?

1. Prepend the sentence with a Start symbol: $\langle s \rangle$
2. Formulate the input sequence as: $(\langle s \rangle, \text{There}, \text{she}, \text{built}, \text{a})$

RNNs: Unrolled step-by-step

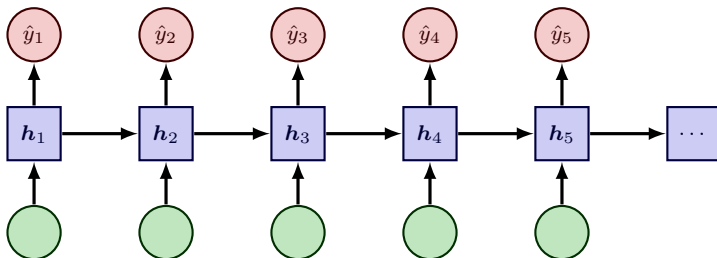
Consider an example of Next Word Prediction:

$$p(\cdot | \text{There she built a })$$

How to build an RNN model for this?

1. Prepend the sentence with a Start symbol: $\langle s \rangle$
2. Formulate the input sequence as: $(\langle s \rangle, \text{There, she, built, a})$

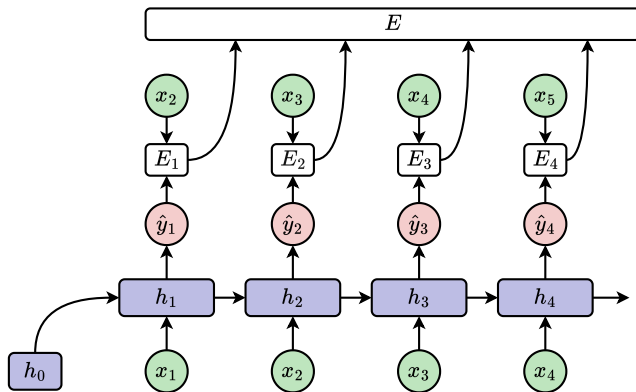
$$h_n = g(V[x_n; h_{n-1}] + c)$$



RNNs: Forward pass

We can run forward equations (black edges) as usual:

$$E = \frac{1}{N} \sum_{n=1} E_n(x_n, \hat{y}_n)$$

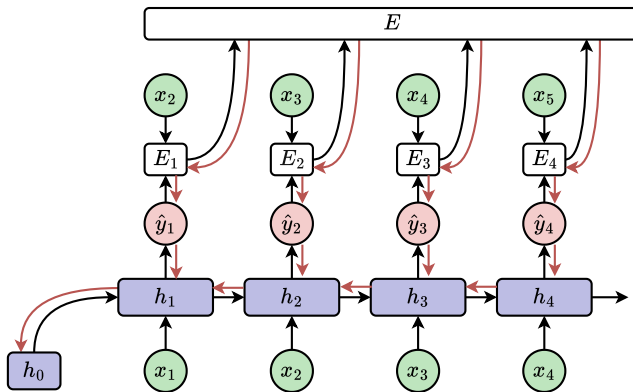


RNNs: Back-Propagation Through Time (1)

We can run forward equations as usual:

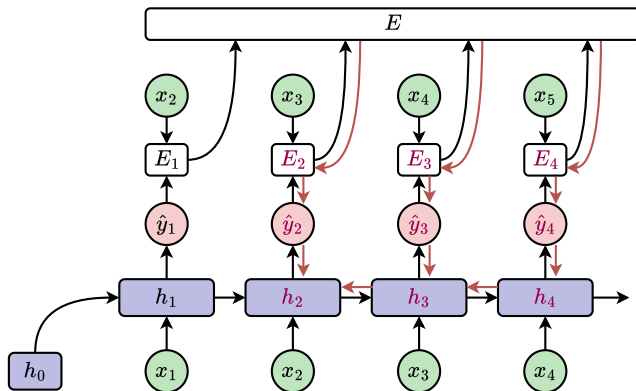
$$E = \frac{1}{N} \sum_{n=1} E_n(x_n, \hat{y}_n)$$

Back-propagation (**red edges**) gives gradients for parameter update:



RNNs: Back-Propagation Through Time (2)

$$\frac{\partial E}{\partial h_2} = \frac{\partial E}{\partial E_2} \frac{\partial E_2}{\partial \hat{y}_2} \frac{\partial \hat{y}_2}{\partial h_2} + \frac{\partial E}{\partial E_3} \frac{\partial E_3}{\partial \hat{y}_3} \frac{\partial \hat{y}_3}{\partial h_3} \frac{\partial h_3}{\partial h_2} + \frac{\partial E}{\partial E_4} \frac{\partial E_4}{\partial \hat{y}_4} \frac{\partial \hat{y}_4}{\partial h_4} \frac{\partial h_4}{\partial h_3} \frac{\partial h_3}{\partial h_2} + \dots$$

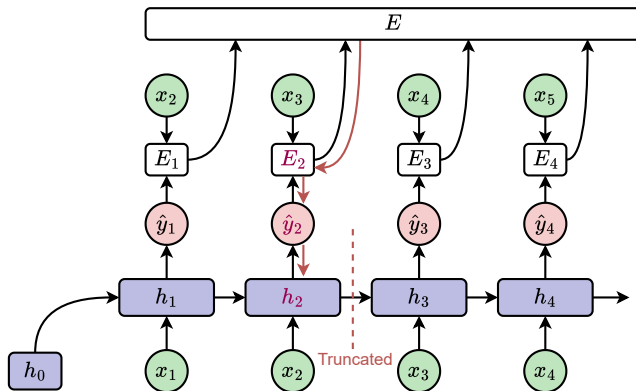


This gradient update is called **Back-Propagation Through Time (BPTT)**: the derivatives at timestep n depend on those at future timesteps:

RNNs: Back-Propagation Through Time (3)

If we break dependencies after a fixed number of timesteps we get **Truncated BPTT**:

$$\frac{\partial E}{\partial \mathbf{h}_2} \approx \frac{\partial E}{\partial E_2} \frac{\partial E_2}{\partial \hat{\mathbf{y}}_2} \frac{\partial \hat{\mathbf{y}}_2}{\partial \mathbf{h}_2}$$



RNNs: Exploding and Vanishing Gradients (1)

An RNN Model needs to discover and represent long-range dependencies:

$p(\text{sandcastle} | \text{Alice went to the beach. There she built a })$

While an RNN model can represent such dependencies in theory, can it learn them?

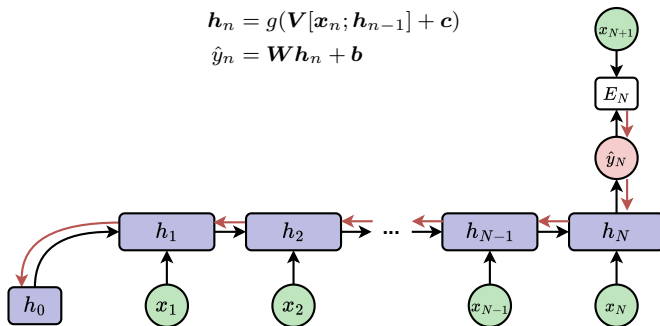
RNNs: Exploding and Vanishing Gradients (1)

An RNN Model needs to discover and represent long-range dependencies:

$p(\text{sandcastle} | \text{Alice went to the beach. There she built a})$

While an RNN model can represent such dependencies in theory, can it learn them?

- Consider the path of partial derivatives linking a change in E_N to changes in h_1 :



$$\frac{\partial E_N}{\partial \mathbf{h}_1} = \frac{\partial E_N}{\partial \hat{\mathbf{y}}_N} \frac{\partial \hat{\mathbf{y}}_N}{\partial \mathbf{h}_N} \left(\prod_{n \in \{N, \dots, 2\}} \frac{\partial \mathbf{h}_n}{\partial \mathbf{h}_{n-1}} \right)$$

RNNs: Exploding and Vanishing Gradients (2)

Assume \mathbf{V} decomposes into \mathbf{V}_x and \mathbf{V}_h and g is an element-wise activation:

$$\mathbf{h}_n = g(\mathbf{V}[\mathbf{x}_n; \mathbf{h}_{n-1}] + \mathbf{c}) = g(\underbrace{\mathbf{V}_x \mathbf{x}_n + \mathbf{V}_h \mathbf{h}_{n-1} + \mathbf{c}}_{\mathbf{z}_n})$$

Calculate the gradients: $\frac{\partial \mathbf{h}_n}{\partial \mathbf{z}_n} = \text{diag}(g'(\mathbf{z}_n))$, $\frac{\partial \mathbf{z}_n}{\partial \mathbf{h}_{n-1}} = \mathbf{V}_h$.

$$\begin{aligned} \frac{\partial E_N}{\partial \mathbf{h}_1} &= \frac{\partial E_N}{\partial \hat{\mathbf{y}}_N} \frac{\partial \hat{\mathbf{y}}_N}{\partial \mathbf{h}_N} \left(\prod_{n \in \{N, \dots, 2\}} \frac{\partial \mathbf{h}_n}{\partial \mathbf{z}_n} \frac{\partial \mathbf{z}_n}{\partial \mathbf{h}_{n-1}} \right) \\ &= \frac{\partial E_N}{\partial \hat{\mathbf{y}}_N} \frac{\partial \hat{\mathbf{y}}_N}{\partial \mathbf{h}_N} \left(\prod_{n \in \{N, \dots, 2\}} \text{diag}(g'(\mathbf{z}_n)) \mathbf{V}_h \right) \end{aligned}$$

For the repeated multiplication of \mathbf{V}_h , if the largest eigenvalue of \mathbf{V}_h is

- ▶ 1, then gradient can be propagated,
- ▶ > 1 , the product will grow exponentially (explode),
- ▶ < 1 , the product shrinks exponentially (vanishes).

The most popular solution: to change the network architecture to include gated units¹

¹Christopher Olah: Understanding LSTM Networks

Outline

Language Models

Recurrent Neural Networks

Transformers

Issues with recurrent neural networks

Lack of parallelizability: Forward and backward passes have $\mathcal{O}(\text{sequence length})$ unparallelizable operations

- ▶ GPUs can perform a bunch of independent computations at once!
- ▶ But future RNN hidden states can't be computed in full before past RNN hidden states have been computed
- ▶ Inhibits training on very large datasets!

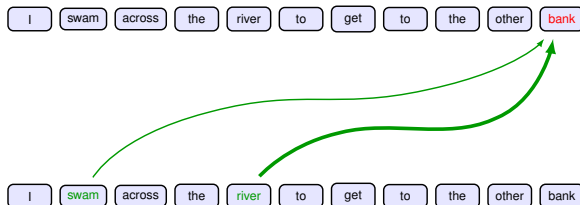
If not recurrence, then what? How about **attention**?

Attention

Consider the following two sentences:

- ▶ I swam across the river to get to the other **bank**.
- ▶ I walked across the road to get cash from the **bank**.

Task: to determine the appropriate interpretation of 'bank'



Attention: a neural network should *attend* to, in other words rely more heavily on, specific words from the rest of the sequence.

Input Representation: word embedding

How to convert the words into a numerical representation?

- ▶ **One-hot encoding:**

- ▶ Define a fixed dictionary and introduce vectors of length equal to dictionary size
- ▶ To encode the k -th word with a vector \mathbf{x}_n having a 1 in position k and 0 elsewhere
- ▶ Results in vectors of very high dimensionality if the dictionary is large

- ▶ **Word embedding:** to map words into a lower-dimensional space

- ▶ Defined a matrix \mathbf{E} of size $D \times K$: D the dimensionality of the embedding space, K the dimensionality of the dictionary

$$\mathbf{v}_n = \mathbf{E}\mathbf{x}_n$$

Vector \mathbf{v}_n is the corresponding column of the matrix \mathbf{E} (can be learned, e.g., *word2vec*²)

- ▶ Learned embedding space often has an even richer semantic structure:

$$\mathbf{v}(\text{Paris}) - \mathbf{v}(\text{France}) \approx \mathbf{v}(\text{Rome}) - \mathbf{v}(\text{Italy})$$

'Paris is to France as Rome is to Italy'

- ▶ As a pre-processing step: part of the overall end-to-end training

²Tomas Mikolov, Kai Chen, Greg Corrado, Jeffrey Dean, Efficient Estimation of Word Representations in Vector Space, 2013

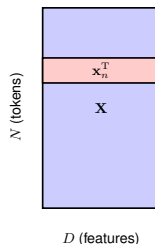
Input Representation

One sentence after embedding:

- ▶ A set of vectors $\{x_n\}$ of dimensionality D , where $n = 1, \dots, N$
- ▶ Vectors are *tokens*, which might correspond to a word or byte pair³

To combine data vectors into a matrix \mathbf{X} of dimensions $N \times D$:

- ▶ n -th row comprises the token vector x_n^\top
- ▶ This matrix represents one set of input tokens.
- ▶ For most applications, we will require a data set containing many sets of tokens.



³Byte pair encoding: https://en.wikipedia.org/wiki/Byte_pair_encoding

Self-attention (1)

Task: to design such input transformation that:

- ▶ Takes \mathbf{X} as input and outputs a transformed matrix \mathbf{Y} of the same dimensionality

$$\mathbf{Y} = \text{TransformerLayer}[\mathbf{X}].$$

- ▶ Can apply multiple times in succession to construct deep networks

Question: Given a set of tokens $\mathbf{x}_1, \dots, \mathbf{x}_N$ in an embedding space, how to map it to another set $\mathbf{y}_1, \dots, \mathbf{y}_N$ having the same token number but in a new embedding space?

- ▶ \mathbf{y}_n should depend not just on \mathbf{x}_n but on $\mathbf{x}_1, \dots, \mathbf{x}_N$, e.g., a linear combination

$$\mathbf{y}_n = \sum_{m=1}^N a_{nm} \mathbf{x}_m,$$

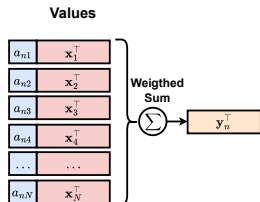
where a_{nm} are called *attention weights*.

- ▶ Coefficients should be close to zero for input tokens that have little influence on the output \mathbf{y}_n and largest for inputs that have most influence.
- ▶ *How to obtain the attention weights?*

Self-attention (2)

Linear combinations:

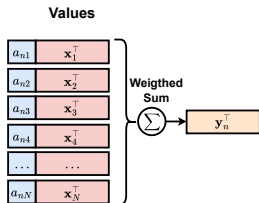
x_n are value vectors to span the output space



Self-attention (2)

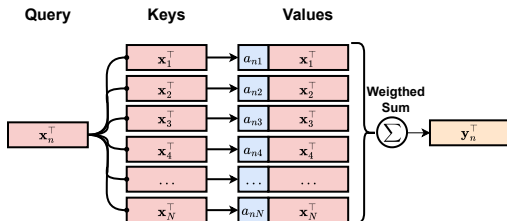
Linear combinations:

x_n are value vectors to span the output space



Transform x_n to y_n :

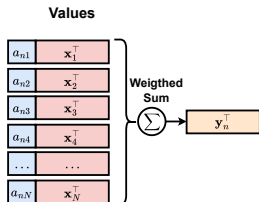
y_n should be closely related to x_n ; use x_n to measure the similarity to all $x_i \forall i$



Self-attention (2)

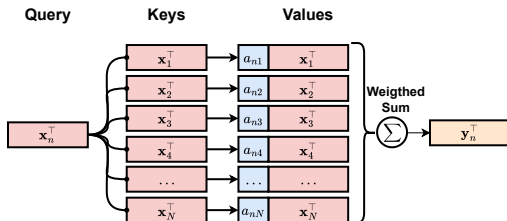
Linear combinations:

x_n are value vectors to span the output space



Transform x_n to y_n :

y_n should be closely related to x_n ; use x_n to measure the similarity to all $x_i \forall i$



Self-attention:

- ▶ Vectors x_n as a *value vector* that will be used to create the output tokens
- ▶ Vectors x_n directly as a *key vector* for input token n
- ▶ Vectors x_m as a *query vector* for output y_m

Namely, vectors x_n are being used as value, key, and query simultaneously.

Self-attention (3)

To measure the similarity between query and key vectors:

- ▶ Dot product $\mathbf{x}_n^\top \mathbf{x}_m$ as similarity measure
- ▶ Define the weighting coefficients a_{nm} by using *softmax* to transform dot products:

$$a_{nm} = \frac{\exp(\mathbf{x}_n^\top \mathbf{x}_m)}{\sum_{m'=1}^N \exp(\mathbf{x}_n^\top \mathbf{x}_{m'})}.$$

Note that in this case there is no probabilistic interpretation of the softmax function and it is simply being used to normalize the attention weights appropriately.

Write in matrix notation by using the data matrix \mathbf{X} and the output matrix \mathbf{Y} :

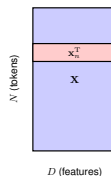
$$\mathbf{Y} = \text{softmax}[\mathbf{X} \mathbf{X}^\top] \mathbf{X}.$$

$\text{softmax}[\mathbf{L}]$ takes the exponential of every element of a matrix \mathbf{L} and then normalises each row independently.

Dot-product self-attention:

- ▶ Same sequence to determine the queries, keys and values
- ▶ Dot product as similarity measure

Transformation from $\{\mathbf{x}_n\}$ to $\{\mathbf{y}_n\}$ is fixed: no learnable parameters.



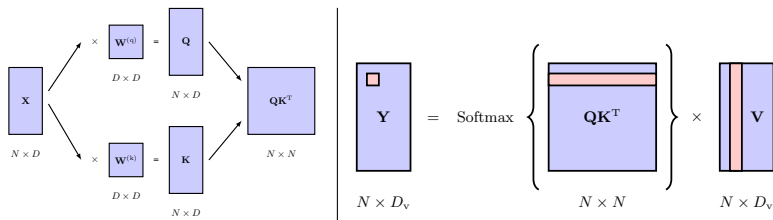
Self-attention with parameters (1)

Defining linear transformations separately for query, key, and value matrices:

$$\mathbf{Q} = \mathbf{X}\mathbf{W}^{(q)}, \quad \mathbf{K} = \mathbf{X}\mathbf{W}^{(k)}, \quad \mathbf{V} = \mathbf{X}\mathbf{W}^{(v)},$$

where weight matrices $\mathbf{W}^{(q)} \in \mathbb{R}^{D \times D_q}$, $\mathbf{W}^{(k)} \in \mathbb{R}^{D \times D_k}$, and $\mathbf{W}^{(v)} \in \mathbb{R}^{D \times D_v}$ represent parameters that will be learned during training. (Typically, $D_k = D$)

$$\mathbf{Y} = \text{softmax}[\underbrace{\mathbf{Q}\mathbf{K}^{\top}}_{N \times N}] \underbrace{\mathbf{V}}_{N \times D_v}$$



We can also include bias parameters in linear transformations (by augmenting the data matrix \mathbf{X} with an additional column and weight matrices with an additional row)

Scaled attention with parameters (2)

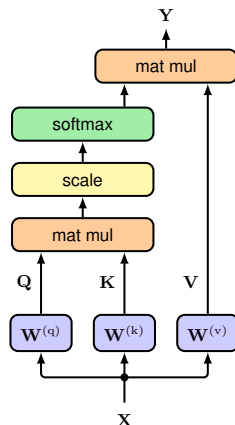
Gradients of softmax function can become exponentially small for inputs of high magnitude, e.g., ReLU activation function.

To prevent this, we can *re-scale the product* before applying the softmax function.

- ▶ If the elements of the query and key vectors are all independent random variables with zero mean and unit variance, the variance of the dot product will be D_k .
- ▶ To normalize the product using the standard deviation given by square root of D_k :

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

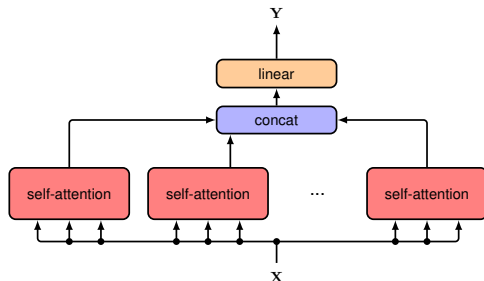
a.k.a **scaled dot-product self-attention**.



Multi-head Attention

- ▶ The attention layer describe so far allows the output vectors to attend to data-dependent patterns of input vectors and is called an *attention head*.
- ▶ Use multiple attention heads in parallel to capture multiple patterns: i.e., identically structured copies with independent learnable parameters, similar to filters in CNN.
- ▶ Suppose we have H heads indexed by $h = 1, \dots, H$:

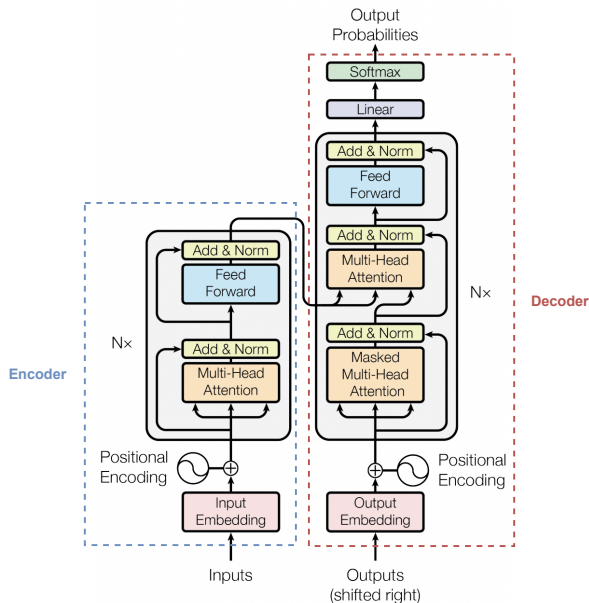
$$H_h = \text{Attention}(Q_h, K_h, V_h).$$



- ▶ The heads are first concatenated into a single matrix and the result is then linearly transformed using a matrix $W^{(o)}$ to give a combined output in the form:

$$Y(X) = \text{Concat}[H_1, \dots, H_H] W^{(o)}.$$

Other layers in transformers (1)



See [paper](#) by Vaswani et al. (2017)

Other layers in transformers (2)

Residual connections and layer normalization:

- ▶ To improve training efficiency, transformers use *residual connections* that bypass the multi-head structure, together with *layer normalization*:

$$\mathbf{Z} = \text{LayerNorm}[\mathbf{Y}(\mathbf{X}) + \mathbf{X}].$$

- ▶ The layer normalization can be replaced by *pre-norm* (more effective):

$$\mathbf{Z} = \mathbf{Y}(\mathbf{X}') + \mathbf{X}, \quad \text{where} \quad \mathbf{X}' = \text{LayerNorm}[\mathbf{X}].$$

Feedforward layers:

- ▶ Non-linearity enters through attention weights via `softmax` function.
- ▶ Output vectors are still constrained to lie in the subspace spanned by the input vectors, which limits the expressive capabilities of the attention layer.
- ▶ To post-process the output of each layer using a standard nonlinear neural network with D inputs and D outputs, denoted $\text{FF}[\cdot]$ (sometimes also as $\text{MLP}[\cdot]$).

$$\tilde{\mathbf{X}} = \text{LayerNorm}[\text{FF}[\mathbf{Z}] + \mathbf{Z}].$$

- ▶ The same shared network is applied to each of the output vectors.
- ▶ The pre-norm form is $\tilde{\mathbf{X}} = \text{FF}[\mathbf{Z}'] + \mathbf{Z}$, where $\mathbf{Z}' = \text{LayerNorm}[\mathbf{Z}]$.

Other layers in transformers (3)

Positional encoding:

- ▶ $W_h^{(q)}$, $W_h^{(k)}$ and $W_h^{(v)}$ are shared across input tokens: permuting the order of input tokens, i.e., rows of X , results in the same permutation of rows of output matrix \tilde{X} .
- ▶ The transformer is *equivariant* with respect to input permutations.
 - ▶ The food was bad, not good at all!
 - ▶ The food was good, not bad at all!

They contain the same tokens but have very different meanings: token order is crucial for most sequential processing tasks.

- ▶ To construct a position encoding vector r_n associated with each input position n :
 - ▶ Concatenation with x_n leads to increased dimensionality.
 - ▶ Simply add the position vectors onto the token vectors to give $\tilde{x}_n = x_n + r_n$.
 - ▶ Residual connection allows it to be passed from one transformer layer to the next

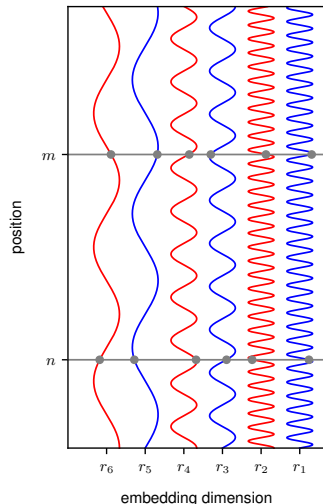
Other layers in transformers (4)

Positional encoding:

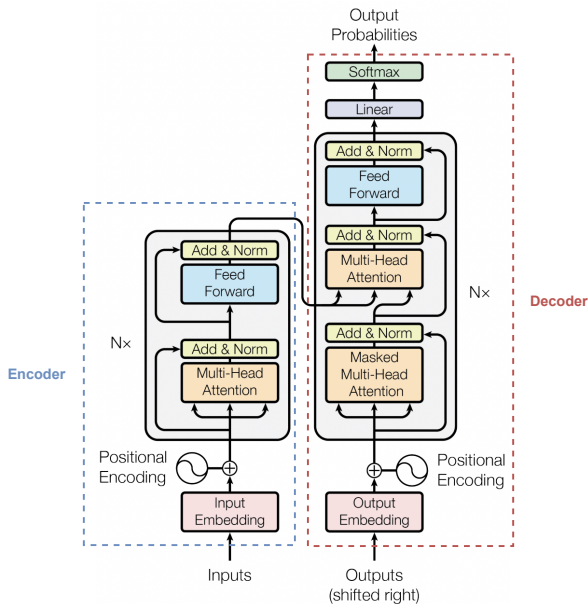
- ▶ Sinusoidal positional encoding (bounded, relative position): for a given position n the associated position-encoding vector has components r_{ni} given by

$$r_{ni} = \begin{cases} \sin\left(\frac{n}{L^{i/D}}\right), & \text{if } i \text{ is even,} \\ \cos\left(\frac{n}{L^{(i-1)/D}}\right), & \text{if } i \text{ is odd.} \end{cases}$$

Encoding at $n + k$ can be represented as a linear combination of encoding at n , in which coefficients do not depend on the absolute position but only on k .



Encoder + Decoder



See [paper](#) by Vaswani et al. (2017)

Encoder transformers (1)

Encoder: takes sequences as input and produces fixed-length vectors.

Bidirectional Encoder Representations from Transformers (BERT):

- ▶ To pre-train a language model using a large corpus of text
- ▶ Then to fine-tune the model using *transfer learning* for a broad range of downstream tasks each of which requires a smaller application-specific dataset.

Training procedure:

- ▶ A randomly chosen subset of the tokens, say 15%, are replaced with a special token denoted $\langle \text{mask} \rangle$. The model is trained to predict the missing tokens at the corresponding output nodes (i.e., *self-supervised learning*).
 - ▶ Original input: I **swam** across the river to get to the **other** bank.
 - ▶ Masked input: I $\langle \text{mask} \rangle$ across the river to get to the $\langle \text{mask} \rangle$ bank.
- ▶ Term 'bidirectional' refers to that the network sees words both before and after the masked word and can use both sources of information to make a prediction.
- ▶ In practice, of the 15% of randomly selected tokens, 80% are replaced with $\langle \text{mask} \rangle$, 10% are replaced with a word selected at random from the vocabulary, and 10% are retained, but they still have to be correctly predicted at the output.⁴

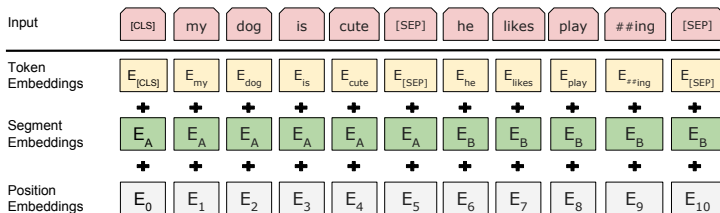
⁴Jacob Devlin, Ming-Wei Chang, Kenton Lee, Kristina Toutanova, BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding, 2018

Encoder transformers (2)

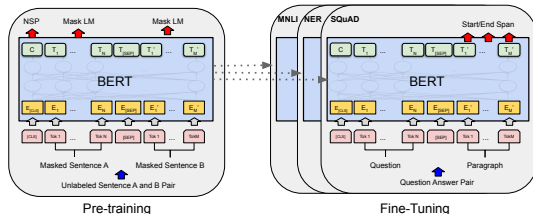
Bidirectional Encoder Representations from Transformers (BERT):

This use of self-supervised learning led to a paradigm shift in which a large model is first pre-trained using unlabelled data and then subsequently fine-tuned using supervised learning based on a much smaller set of labelled data.

► Input representation:

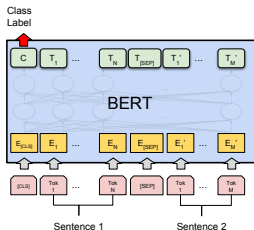


► Pre-training and fine-tuning:

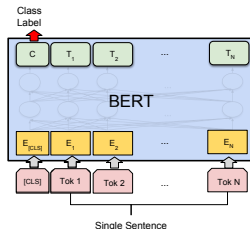


Encoder transformers (3)

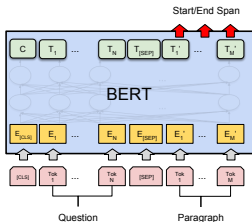
Bidirectional Encoder Representations from Transformers (BERT):



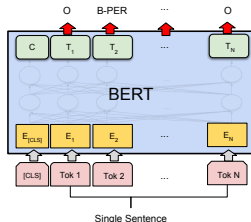
(a) Sentence Pair Classification Tasks:
MNLI, QQP, QNLI, STS-B, MRPC,
RTE, SWAG



(b) Single Sentence Classification Tasks:
SST-2, CoLA



(c) Question Answering Tasks:
SQuAD v1.1

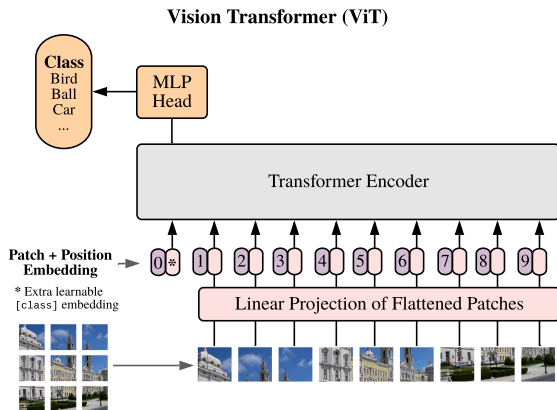


(d) Single Sentence Tagging Tasks:
CoNLL-2003 NER

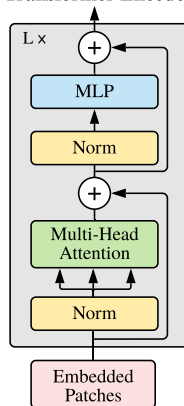
Encoder transformers (4)

Vision transformers

- ▶ Slices an image to patches
- ▶ Embed all patches to form tokens
- ▶ Apply encoder transformers with pre-training



Transformer Encoder



Decoder transformers (1)

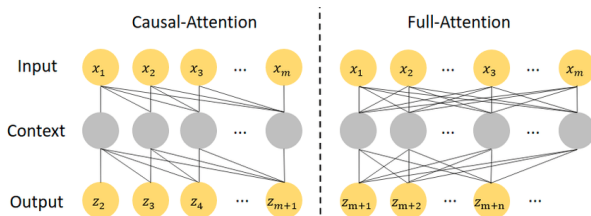
Decoder: used as generative models that create output sequence of tokens,

Autoregressive model

- ▶ Take as input a sequence consisting of the first $n - 1$ tokens, and its corresponding output represents the conditional distribution for token n
- ▶ Draw a sample from this distribution and extend the sequence to n tokens
- ▶ Feed back the new sequence through model to give a distribution over token $n + 1$

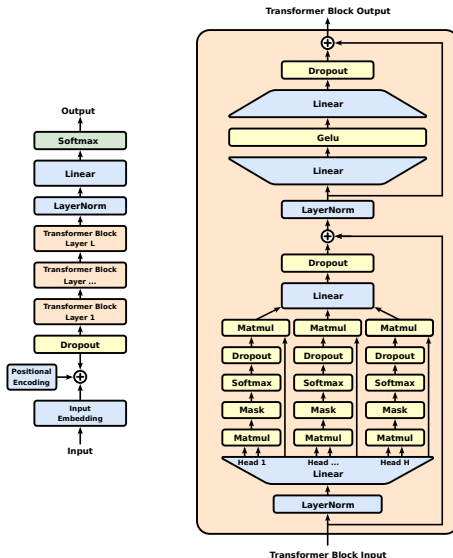
Generative Pre-trained Transformers (GPT)

- ▶ Shift input sequence to the right by one step: x_n corresponds to y_{n+1} , with target x_{n+1} , and prepend a special token $\langle \text{start} \rangle$ in the first position of input sequence
- ▶ *Causal attention*: set to zero all of attention weights of any later token in the sequence, i.e., setting corresponding pre-activation values to $-\infty$ before softmax



Decoder transformers (2)

Generative Pre-trained Transformer (GPT)



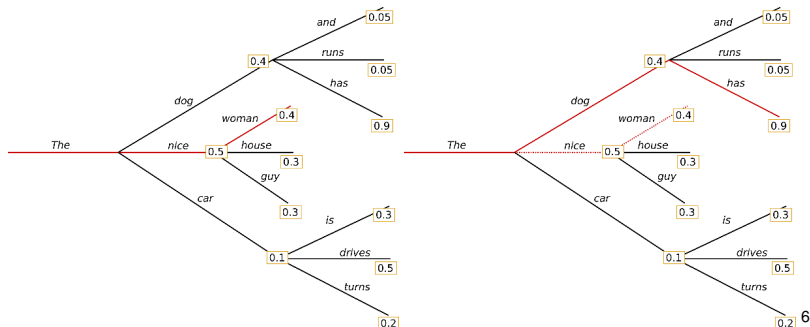
5

Decoder transformers (3)

How to generate the next token?

Sampling scheme:

- ▶ *Greedy search*: simply to select tokens with the *highest probability*.
- ▶ *Beam search*: to maintain a set of B hypotheses, where B is called the *beam width*, each consisting of a sequence up to step $n + B$. To select tokens with the *highest total probability* of the textended sequences.



6

Sequence-to-sequence (seq2seq) transformers

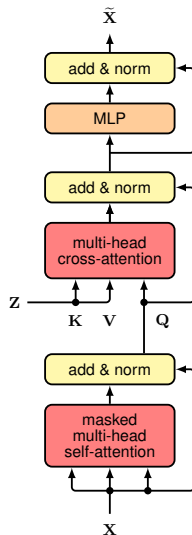
Cross-attention: combine encoder with decoder

Consider sequence-to-sequence modeling (Seq2seq):
For example, *translation task*: translating an English sentence into a Dutch sentence.

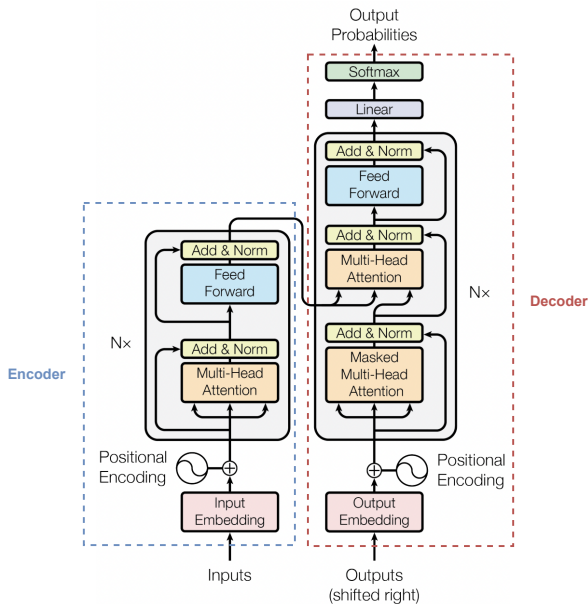
1. Use an encoder transformer to map the input token sequence into a suitable internal representation Z
2. Use a decoder transformer to generate the token sequence corresponding to the Dutch output, token by token
3. Condition the output on the entire input sequence corresponding to the English sentence

To incorporate Z into the decoder, we use the **cross attention** in *decoder*:

- ▶ Query vectors come from the sequence being generated, i.e., from decoder
- ▶ Key and value vectors come from the sequence represented by Z , i.e., from encoder



Full transformers



See [paper](#) by Vaswani et al. (2017)

Recap

Language Models

Recurrent Neural Networks

Transformers

END LECTURE