



UPPSALA
UNIVERSITET

Tranformers and BERT

Måns Magnusson
Statistiska Institutionen
Uppsala Universitet

17th of June 2024

- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa



UPPSALA
UNIVERSITET

- **Introduction**
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

Section 1

Introduction



● Introduction

● (Recap) Feed-Forward Neural Networks

- Hidden Units
- Architecture design

● The Transformer

- Attention
- Multi-Head Attention
- Positional encoding
- Add and Normalize

● Transformer-Encoder Models

- BERT
- RoBERTa

The Development of NLP

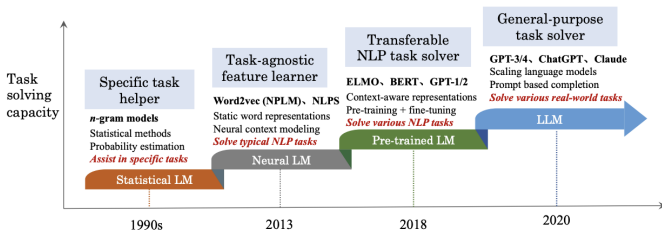


Figure: The development timeline (Zhao et al., 2023)



UPPSALA
UNIVERSITET

Why and when neural networks?

- Introduction

- (Recap) Feed-Forward Neural Networks

- Hidden Units
- Architecture design

- The Transformer

- Attention
- Multi-Head Attention
- Positional encoding
- Add and Normalize

- Transformer-Encoder Models

- BERT
- RoBERTa



UPPSALA
UNIVERSITET

Why and when neural networks?

- Introduction

- (Recap) Feed-Forward Neural Networks

- Hidden Units
- Architecture design

- The Transformer

- Attention
- Multi-Head Attention
- Positional encoding
- Add and Normalize

- Transformer-Encoder Models

- BERT
- RoBERTa

- Learning feature representations



UPPSALA
UNIVERSITET

Why and when neural networks?

- Introduction

- (Recap) Feed-Forward Neural Networks

- Hidden Units
- Architecture design

- The Transformer

- Attention
- Multi-Head Attention
- Positional encoding
- Add and Normalize

- Transformer-Encoder Models

- BERT
- RoBERTa

- Learning **feature representations**
- Good for "sensor" data



UPPSALA
UNIVERSITET

Why and when neural networks?

- Introduction

- (Recap) Feed-Forward Neural Networks

- Hidden Units
- Architecture design

- The Transformer

- Attention
- Multi-Head Attention
- Positional encoding
- Add and Normalize

- Transformer-Encoder Models

- BERT
- RoBERTa

- Learning **feature representations**
- Good for "sensor" data
- Needs a lot of data to learn **complex representations** (image, text, audio)



- Introduction

- (Recap) Feed-Forward Neural Networks

- Hidden Units
- Architecture design

- The Transformer

- Attention
- Multi-Head Attention
- Positional encoding
- Add and Normalize

- Transformer-Encoder Models

- BERT
- RoBERTa

Learning Representations

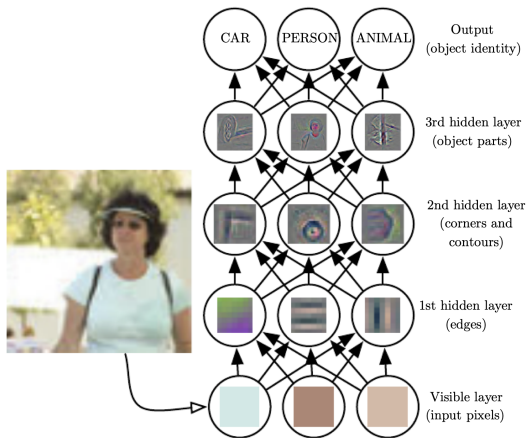


Figure: Learning representations can be crucial (Goodfellow et al, 2017, Fig. 1.2)



UPPSALA
UNIVERSITET

Different Network Architectures

- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa
- Different networks for different purposes
 - **Feed-Forward Neural Network**: Basic building block



UPPSALA
UNIVERSITET

Different Network Architectures

- Introduction

- (Recap) Feed-Forward Neural Networks

- Hidden Units
- Architecture design

- The Transformer

- Attention
- Multi-Head Attention
- Positional encoding
- Add and Normalize

- Transformer-Encoder Models

- BERT
- RoBERTa

- Different networks for different purposes

- **Feed-Forward Neural Network:** Basic building block
- **Convolutional Neural Networks:** Computer Vision



- Introduction

- (Recap) Feed-Forward Neural Networks

- Hidden Units
- Architecture design

- The Transformer

- Attention
- Multi-Head Attention
- Positional encoding
- Add and Normalize

- Transformer-Encoder Models

- BERT
- RoBERTa

- Different networks for different purposes

- **Feed-Forward Neural Network:** Basic building block
- **Convolutional Neural Networks:** Computer Vision
- **Transformers:** Textual data



UPPSALA
UNIVERSITET

- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

Section 2

(Recap) Feed-Forward Neural Networks



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

The Feed-Forward Network

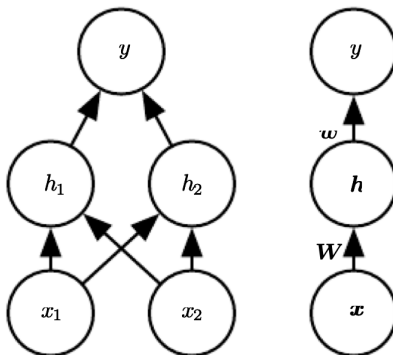


Figure: A simple feed-forward network (Goodfellow et al, 2017, Fig. 6.2)

Important concepts:

Layers, neurons, input, output, weights, bias, architecture



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

The Feed-Forward Network

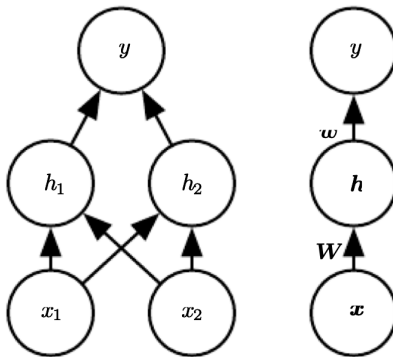


Figure: A simple feed-forward network (Goodfellow et al, 2017, Fig. 6.2)

In mathematical notation:

$$y_i = \mathbf{w}^T g(\mathbf{W}^T \mathbf{x}_i + \mathbf{b}_1) + \mathbf{b}_2,$$

where \mathbf{w} , \mathbf{W} , \mathbf{b}_1 and \mathbf{b}_2 is learned/estimated



UPPSALA
UNIVERSITET

The Feed-Forward Network

$$y_i = \mathbf{w}^T g(\mathbf{W}^T \mathbf{x}_i + \mathbf{b}_1) + \mathbf{b}_2$$

- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa



$$y_i = \mathbf{w}^T g(\mathbf{W}^T \mathbf{x}_i + \mathbf{b}_1) + \mathbf{b}_2$$

$$W = \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}, w = \begin{pmatrix} 1 \\ -2 \end{pmatrix}, b_1 = \begin{pmatrix} 1 \\ -1 \end{pmatrix}, b_2 = (0)$$

- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa



$$y_i = \mathbf{w}^T g(\mathbf{W}^T \mathbf{x}_i + \mathbf{b}_1) + \mathbf{b}_2$$

$$W = \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}, w = \begin{pmatrix} 1 \\ -2 \end{pmatrix}, b_1 = \begin{pmatrix} 1 \\ -1 \end{pmatrix}, b_2 = (0)$$

$$g(z) = \text{ReLU}(z) = \max(0, z)$$

- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa



$$y_i = \mathbf{w}^T g(\mathbf{W}^T \mathbf{x}_i + \mathbf{b}_1) + \mathbf{b}_2$$

$$W = \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}, w = \begin{pmatrix} 1 \\ -2 \end{pmatrix}, b_1 = \begin{pmatrix} 1 \\ -1 \end{pmatrix}, b_2 = (0)$$

$$g(z) = \text{ReLU}(z) = \max(0, z)$$

$$x_i = \begin{pmatrix} 0 \\ 0 \end{pmatrix},$$

- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa



$$y_i = \mathbf{w}^T g(\mathbf{W}^T \mathbf{x}_i + \mathbf{b}_1) + \mathbf{b}_2$$

$$W = \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}, w = \begin{pmatrix} 1 \\ -2 \end{pmatrix}, b_1 = \begin{pmatrix} 1 \\ -1 \end{pmatrix}, b_2 = (0)$$

$$g(z) = \text{ReLU}(z) = \max(0, z)$$

$$x_i = \begin{pmatrix} 0 \\ 0 \end{pmatrix},$$

$$y_i = \begin{pmatrix} 1 \\ -2 \end{pmatrix}^T \text{ReLU} \left[\begin{pmatrix} 0 \\ 0 \end{pmatrix} + \begin{pmatrix} 1 \\ -1 \end{pmatrix} \right] + (0)$$

- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa



$$y_i = \mathbf{w}^T g(\mathbf{W}^T \mathbf{x}_i + \mathbf{b}_1) + \mathbf{b}_2$$

$$W = \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}, w = \begin{pmatrix} 1 \\ -2 \end{pmatrix}, b_1 = \begin{pmatrix} 1 \\ -1 \end{pmatrix}, b_2 = (0)$$

$$g(z) = \text{ReLU}(z) = \max(0, z)$$

$$x_i = \begin{pmatrix} 0 \\ 0 \end{pmatrix},$$

$$y_i = \begin{pmatrix} 1 \\ -2 \end{pmatrix}^T \text{ReLU} \left[\begin{pmatrix} 0 \\ 0 \end{pmatrix} + \begin{pmatrix} 1 \\ -1 \end{pmatrix} \right] + (0)$$

$$y_i = \begin{pmatrix} 1 \\ -2 \end{pmatrix}^T \begin{pmatrix} 1 \\ 0 \end{pmatrix} + (0) = 1$$

- Introduction

- (Recap) Feed-Forward Neural Networks

- Hidden Units
- Architecture design

- The Transformer

- Attention
- Multi-Head Attention
- Positional encoding
- Add and Normalize

- Transformer-Encoder Models

- BERT
- RoBERTa



Output units (g_L)

- Depend on the data y

- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa



Output units (g_L)

- Depend on the data y
- Linear units for regression

$$\hat{y} = \mathbf{wh} + \mathbf{b}$$

- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

Output units (g_L)

- Depend on the data y
- Linear units for regression

$$\hat{y} = \mathbf{wh} + \mathbf{b}$$

- Bernoulli units for binary classification

$$\hat{y} = \sigma(\mathbf{wh} + \mathbf{b}),$$

where

$$\sigma(z) = \frac{1}{1 + \exp(-z)},$$

i.e. the **logistic** or **sigmoid** function.

- Other likelihoods can be used, such as Multinomial, Poisson, etc.



Activation functions (g_l)

- Historically $g(z)$ has been the **sigmoid** or or hyperbolic tangent (\tanh)

$$g_{\text{sigmoid}}(z) = \frac{e^z}{e^z + 1} = \frac{1}{1 + e^{-z}}$$

$$g_{\text{tanh}}(z) = \frac{\sinh z}{\cosh z} = \frac{e^{2z} - 1}{e^{2z} + 1}$$

- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa



Activation functions (g_l)

- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

- Historically $g(z)$ has been the **sigmoid** or or hyperbolic tangent (\tanh)

$$g_{\text{sigmoid}}(z) = \frac{e^z}{e^z + 1} = \frac{1}{1 + e^{-z}}$$

$$g_{\text{tanh}}(z) = \frac{\sinh z}{\cosh z} = \frac{e^{2z} - 1}{e^{2z} + 1}$$

- Today, variants of Rectified linear unit (ReLU) is common
 - Easier to estimate with SGD
 - Easier for deep models



Activation functions (g_l): ReLU

$$g_{\text{ReLU}}(z) = \max(0, z)$$

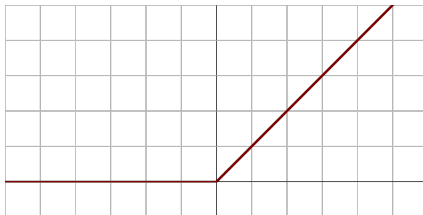


Figure: Rectified Linear Unit (Wikipedia)

- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa



Activation functions (g_l): GeLU

$$g_{\text{GeLU}}(z) = z\Phi(z)$$

where $\Phi(z)$ is a standard Gaussian.

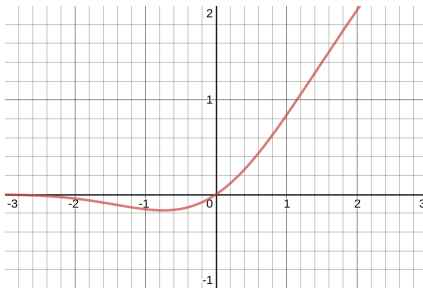


Figure: Gaussian Error Linear Unit (Wikipedia)

- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

Activation functions (g_l): ELU

$$g_{\text{ELU}}(z) = \begin{cases} \alpha (e^z - 1) & \text{if } z \leq 0 \\ z & \text{if } z > 0 \end{cases}$$

where α is commonly 1.

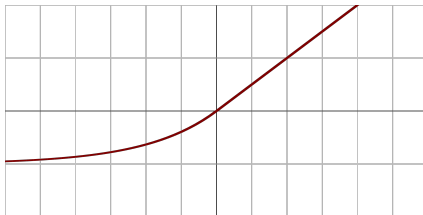


Figure: Exponential Linear Unit (Wikipedia)



UPPSALA
UNIVERSITET

Architecture design

- Introduction
 - (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
 - The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
 - Transformer-Encoder Models
 - BERT
 - RoBERTa
- Architecture: the overall structure of the network
 - Choices:
 - How many layers?
 - How many hidden units in each layer?
 - Activation functions?



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

Depth matters

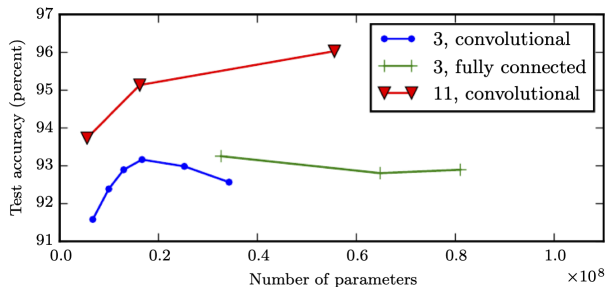


Figure: Depth vs. no of parameters (Goodfellow et al, 2017, Fig. 6.3)



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

Depth matters II

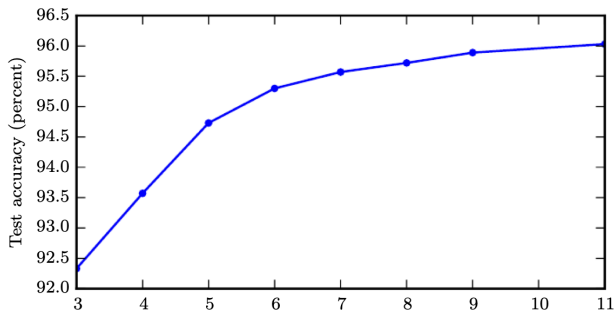


Figure: Effect of depth on accuracy (Goodfellow et al, 2017, Fig. 6.6)



UPPSALA
UNIVERSITET

- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- **The Transformer**
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

Section 3

The Transformer



UPPSALA
UNIVERSITET

The Transformer

- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- **The Transformer**
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

- Introduced by Vaswani et al. (2017):
Attention is all you need.



UPPSALA
UNIVERSITET

The Transformer

- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- **The Transformer**
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

- Introduced by Vaswani et al. (2017):
Attention is all you need.
- Behind the recent progress in NLP: BERT, Llama, GPT, etc.



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- **The Transformer**
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

- Introduced by Vaswani et al. (2017):
Attention is all you need.
- Behind the recent progress in NLP: BERT, Llama, GPT, etc.
- Benefits for textual data:
 - Enables more GPU **parallelism**
 - Better handling of **long-range dependencies**
 - Enable **transfer learning** for text data
 - Enables **deeper** networks for text



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- **The Transformer**
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

- Introduced by Vaswani et al. (2017):
Attention is all you need.
- Behind the recent progress in NLP: BERT, Llama, GPT, etc.
- Benefits for textual data:
 - Enables more GPU **parallelism**
 - Better handling of **long-range dependencies**
 - Enable **transfer learning** for text data
 - Enables **deeper** networks for text
- I will rely heavily on images from Allamar (2018) **The Illustrated Transformer** (recommended)



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- **The Transformer**
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

A Sequence-to-Sequence Model



Figure: The basic block (Allamar, 2018)



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

Stacked Encoder-Decoder Structure

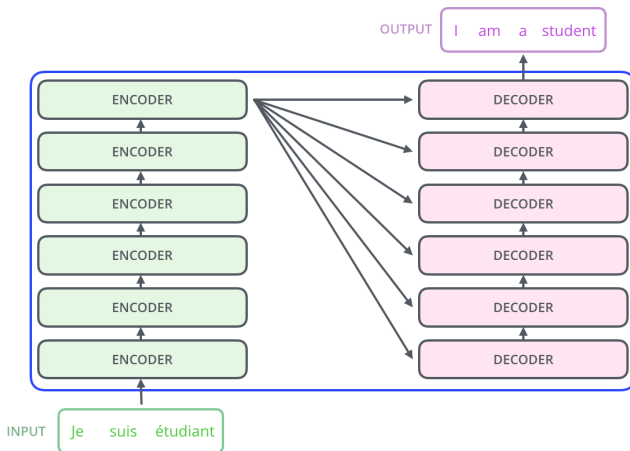


Figure: The Transformer layers (Allamar, 2018)



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

Transformer

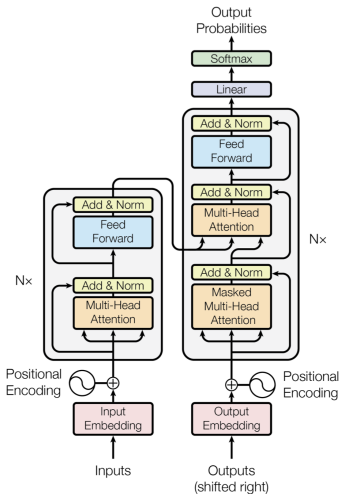


Figure: The Transformer Architecture (Vaswani et al., 2017)



The encoder vs. the decoder

- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- **The Transformer**
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

- Encoder:
 - Input: words
 - Output: contextualized embeddings
- Decoder:
 - Input: **previous words** (and contextualized embeddings from encoder)
 - Output: next word prediction



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

The Transformer Layer (Encoder layer)

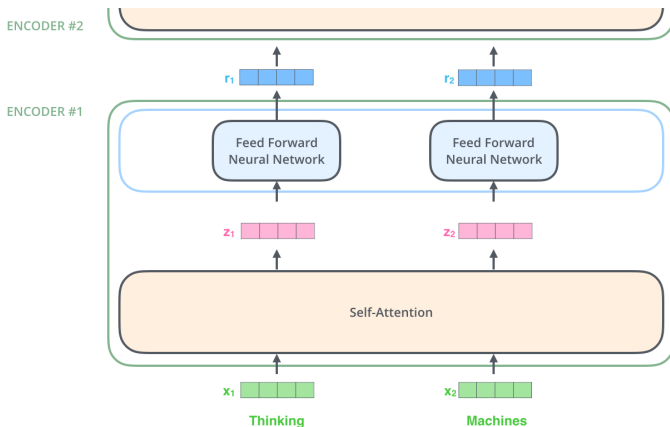


Figure: The Encoder Layer (Alammar, 2018b)



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

Scaled Dot-Product Attention

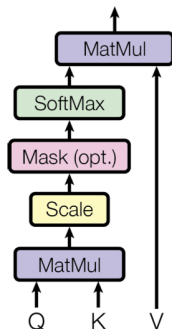


Figure: Scaled Dot-Product Attention (Vaswani et al., 2017)

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



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - **Attention**
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

- (Q)uery: Word i query other words
- (K)ey: The other words return their key to i
- (V)alue: The value of the other words to i



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

Computing Q, V and K

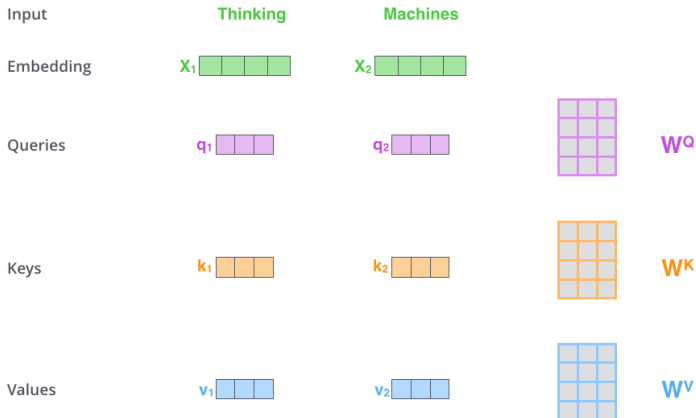


Figure: Attention heads (Alammar, 2018b)



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

Computing Self-Attention

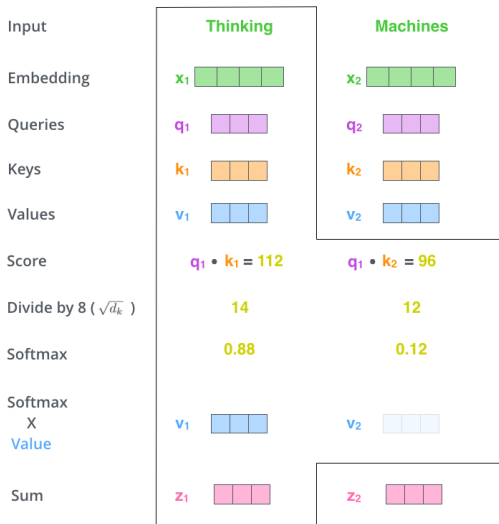


Figure: Attention (Alammar, 2018b)



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - **Multi-Head Attention**
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

Multi-Head Attention

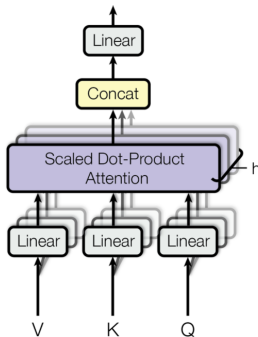


Figure: Scaled Dot-Product Attention (Vaswani et al., 2017)



Attentions Heads

- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

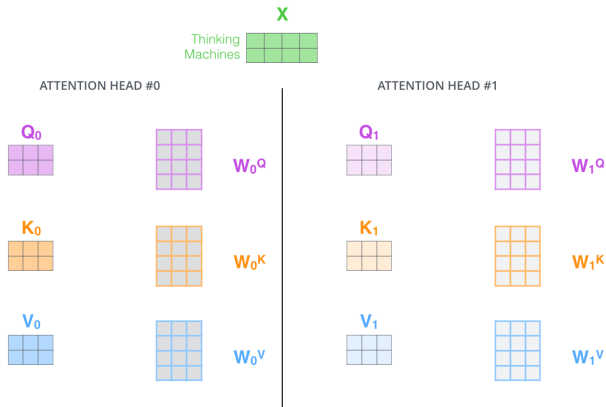


Figure: Attention heads (Alammar, 2018b)



- Introduction

- (Recap) Feed-Forward Neural Networks

- Hidden Units
- Architecture design

- The Transformer

- Attention
- Multi-Head Attention
- Positional encoding
- Add and Normalize

- Transformer-Encoder Models

- BERT
- RoBERTa

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

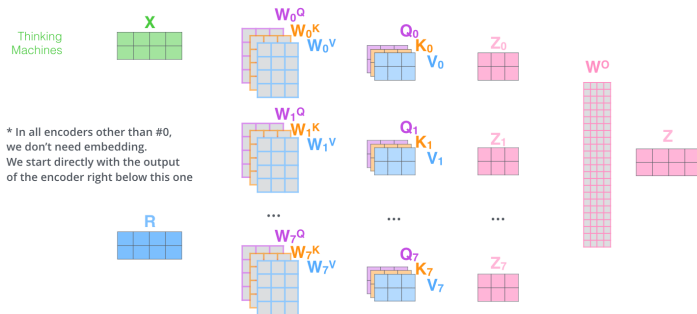


Figure: Attention heads (Alammar, 2018b)



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

Multi-Head Attention example

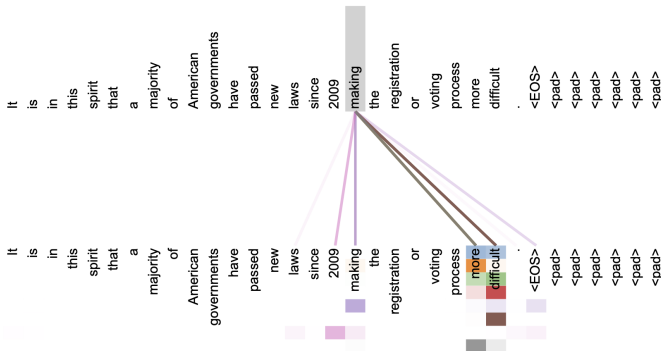


Figure 3: An example of the attention mechanism following long-distance dependencies in the encoder self-attention in layer 5 of 6. Many of the attention heads attend to a distant dependency of the verb 'making', completing the phrase 'making...more difficult'. Attentions here shown only for the word 'making'. Different colors represent different heads. Best viewed in color.

Figure: Attention (Vaswani et al., 2017)



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - **Positional encoding**
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

Positional Encoding

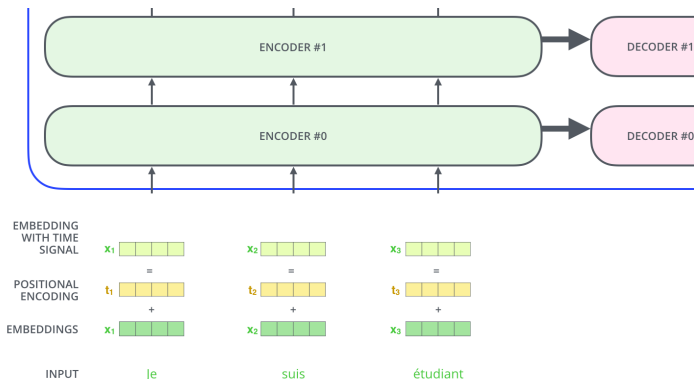


Figure: Attention heads (Alammar, 2018b)



(Absolute) Positional Encoding

"The boy hit the ball" vs "The ball hit the boy"



Figure: Adding positional encodings to embeddings (Alammar, 2018b)



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - **Positional encoding**
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

"The boy hit the ball" vs "The ball hit the boy"

1. Absolute position encoding
2. Relative position encoding:
The distance between tokens are added as bias terms



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - **Positional encoding**
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

"The boy hit the ball" vs "The ball hit the boy"

1. Absolute position encoding
2. Relative position encoding:
The distance between tokens are added as bias terms
3. Rotational positional encoding (RoPE, Su et al., 2022):
Q and K are rotated by based on distance



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

Add and Normalize

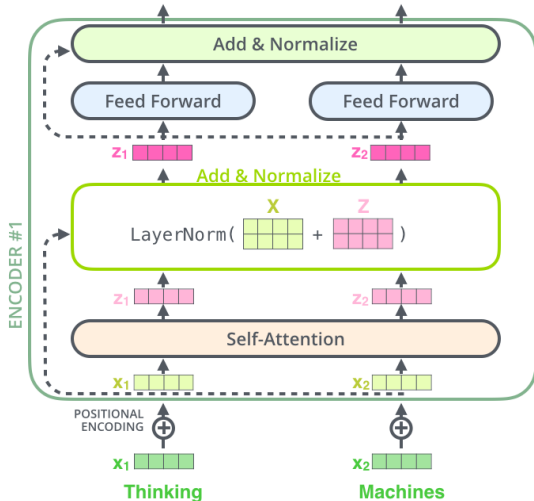


Figure: Add and Normalize (Alammar, 2018b)



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

Transformer

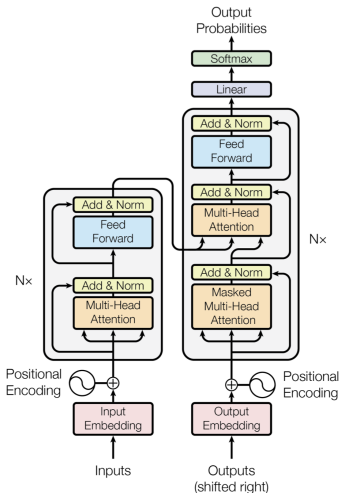


Figure: The Transformer Architecture (Vaswani et al., 2017)



UPPSALA
UNIVERSITET

Tokenization

- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

- Subword tokenization is commonly used



- Introduction
 - (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
 - The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
 - Transformer-Encoder Models
 - BERT
 - RoBERTa
- Subword tokenization is commonly used
 - The main problem with tokenization
 1. Very large vocabulary size
 2. Out-of-vocabulary (OOV) tokens
 3. Different meanings of very similar words



- Introduction
 - (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
 - The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
 - Transformer-Encoder Models
 - BERT
 - RoBERTa
- Subword tokenization is commonly used
 - The main problem with tokenization
 1. Very large vocabulary size
 2. Out-of-vocabulary (OOV) tokens
 3. Different meanings of very similar words
 - Two common approaches:
 1. Byte-pair encoding (GPT-2, RoBERTa)
 2. WordPiece (BERT)



UPPSALA
UNIVERSITET

Byte-pair encoding

- Gage, Philip (1994). "A New Algorithm for Data Compression"

- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa



- Gage, Philip (1994). "A New Algorithm for Data Compression"
- Encode the most common pairs iteratively
 1. look for the most frequent pairing
 2. merge them
 3. repeat (until token or iteration limit)

- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

Byte-pair encoding

- Gage, Philip (1994). "A New Algorithm for Data Compression"
- Encode the most common pairs iteratively
 1. look for the most frequent pairing
 2. merge them
 3. repeat (until token or iteration limit)
- Example (Wikipedia): aaabdaaabc



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

Byte-pair encoding

- Gage, Philip (1994). "A New Algorithm for Data Compression"
- Encode the most common pairs iteratively
 1. look for the most frequent pairing
 2. merge them
 3. repeat (until token or iteration limit)
- Example (Wikipedia): aaabdaaabac

Step 1: ZabdZabac

Z=aa



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

Byte-pair encoding

- Gage, Philip (1994). "A New Algorithm for Data Compression"
- Encode the most common pairs iteratively
 1. look for the most frequent pairing
 2. merge them
 3. repeat (until token or iteration limit)
- Example (Wikipedia): aaabdaaabac

Step 1: ZabdZabac

Z=aa

Step 2: ZYdZYac

Z=aa

Y=ab



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

Byte-pair encoding

- Gage, Philip (1994). "A New Algorithm for Data Compression"
- Encode the most common pairs iteratively
 1. look for the most frequent pairing
 2. merge them
 3. repeat (until token or iteration limit)
- Example (Wikipedia): aaabdaaabac

Step 1: ZabdZabac

Z=aa

Step 2: ZYdZYac

Z=aa

Y=ab

Step 3: XdXac

Z=aa

Y=ab

X=ZY



UPPSALA UNIVERSITET

- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

Byte-pair encoding

1. look for the most frequent pairing
 2. merge them
 3. repeat (until token or iteration limit)
- Example : 9:text_, 10:texting_,11:context_



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

Byte-pair encoding

1. look for the most frequent pairing
 2. merge them
 3. repeat (until token or iteration limit)
- Example : 9:text_, 10:texting_,11:context_
Step 1 (most common: "te"): {te}



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

Byte-pair encoding

1. look for the most frequent pairing
 2. merge them
 3. repeat (until token or iteration limit)
- Example : 9:text_, 10:texting_,11:context_

Step 1 (most common: "te"): {te}

...

Step i (most common: "text_"): {text_}



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

Byte-pair encoding

1. look for the most frequent pairing
 2. merge them
 3. repeat (until token or iteration limit)
- Example : 9:text_, 10:texting_,11:context_

Step 1 (most common: "te"): {te}

...

Step i (most common: "text_"): {text_}

...

Step j (most common: "con"): {text_,con}



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

Byte-pair encoding

1. look for the most frequent pairing
2. merge them
3. repeat (until token or iteration limit)

- Example : 9:text_, 10:texting_,11:context_

Step 1 (most common: "te"): {te}

...

Step i (most common: "text_"): {text_}

...

Step j (most common: "con"): {text_,con}

...

Step k (most common: "texting"): {text_,con,texting}



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

Byte-pair encoding

1. look for the most frequent pairing
2. merge them
3. repeat (until token or iteration limit)

- Example : 9:text_, 10:texting_,11:context_

Step 1 (most common: "te"): {te}

...

Step i (most common: "text_"): {text_}

...

Step j (most common: "con"): {text_,con}

...

Step k (most common: "texting"): {text_,con,texting}



UPPSALA UNIVERSITET

WordPiece

- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa
- BPE difficulty: Which pair to choose (if they are approx. equally frequent)?



UPPSALA
UNIVERSITET

WordPiece

- Introduction
 - (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
 - The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
 - Transformer-Encoder Models
 - BERT
 - RoBERTa
- BPE difficulty: Which pair to choose (if they are approx. equally frequent)?
 - Schuster and Kaisuke (2012) present the WordPiece model



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

- BPE difficulty: Which pair to choose (if they are approx. equally frequent)?
- Schuster and Kaisuke (2012) present the WordPiece model
- Let $P(i, j)$ be the probability of observing the pair ij and $P(i)$ observing i .



- Introduction
 - (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
 - The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
 - Transformer-Encoder Models
 - BERT
 - RoBERTa
- BPE difficulty: Which pair to choose (if they are approx. equally frequent)?
 - Schuster and Kaisuke (2012) present the WordPiece model
 - Let $P(i, j)$ be the probability of observing the pair ij and $P(i)$ observing i .
 - **BPE**: Choose highest $P(i, j)$



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

- BPE difficulty: Which pair to choose (if they are approx. equally frequent)?
- Schuster and Kaisuke (2012) present the WordPiece model
- Let $P(i, j)$ be the probability of observing the pair ij and $P(i)$ observing i .
- **BPE**: Choose highest $P(i, j)$
- **Wordpiece**: Choose highest $P(i, j)/(P(i)P(j))$



UPPSALA
UNIVERSITET

- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

Section 4

Transformer-Encoder Models



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

Transformer-Encoder Models

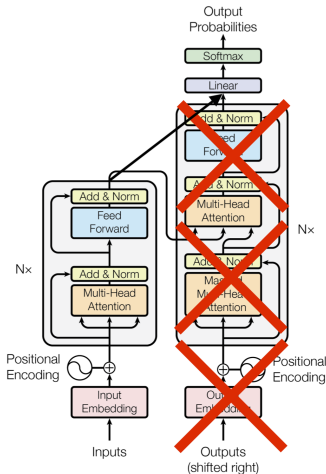


Figure: The Transformer Architecture (Vaswani et al., 2017)



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

Transformer-Encoder Models

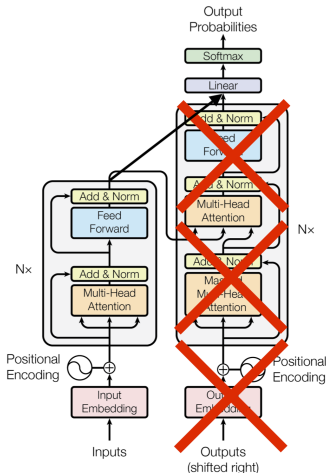


Figure: The Transformer Architecture (Vaswani et al., 2017)

- Common models are BERT and RoBERTa



UPPSALA UNIVERSITET

BERT

- Bidirectional Encoder Representations from Transformers (BERT)
- Introduced in Devlin et al. (2018)

- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

- Bidirectional Encoder Representations from Transformers (BERT)
- Introduced in Devlin et al. (2018)
- **State-of-the-Art** in many text prediction tasks, such as
 - Named-Entity Recognition
 - Text Classification



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

- Bidirectional Encoder Representations from Transformers (BERT)
- Introduced in Devlin et al. (2018)
- **State-of-the-Art** in many text prediction tasks, such as
 - Named-Entity Recognition
 - Text Classification
- Many flavors, such as RoBERTa, ALBERT, etc.



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

- Bidirectional Encoder Representations from Transformers (BERT)
- Introduced in Devlin et al. (2018)
- **State-of-the-Art** in many text prediction tasks, such as
 - Named-Entity Recognition
 - Text Classification
- Many flavors, such as RoBERTa, ALBERT, etc.
- **Pre-trained** on a large corpus



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

- Bidirectional Encoder Representations from Transformers (BERT)
- Introduced in Devlin et al. (2018)
- **State-of-the-Art** in many text prediction tasks, such as
 - Named-Entity Recognition
 - Text Classification
- Many flavors, such as RoBERTa, ALBERT, etc.
- **Pre-trained** on a large corpus
- Then **fine-tuned** for a specific problem



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

- Bidirectional Encoder Representations from Transformers (BERT)
- Introduced in Devlin et al. (2018)
- **State-of-the-Art** in many text prediction tasks, such as
 - Named-Entity Recognition
 - Text Classification
- Many flavors, such as RoBERTa, ALBERT, etc.
- **Pre-trained** on a large corpus
- Then **fine-tuned** for a specific problem
- Available English, Swedish and many other languages (The National Library)



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

- Bidirectional Encoder Representations from Transformers (BERT)
- Introduced in Devlin et al. (2018)
- **State-of-the-Art** in many text prediction tasks, such as
 - Named-Entity Recognition
 - Text Classification
- Many flavors, such as RoBERTa, ALBERT, etc.
- **Pre-trained** on a large corpus
- Then **fine-tuned** for a specific problem
- Available English, Swedish and many other languages (The National Library)
- And again, I rely on Alammari (2018) **The illustrated BERT**

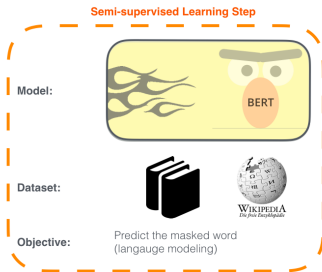


- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

BERT and transfer learning

1 - **Semi-supervised** training on large amounts of text (books, wikipedia..etc).

The model is trained on a certain task that enables it to grasp patterns in language. By the end of the training process, BERT has language-processing abilities capable of empowering many models we later need to build and train in a supervised way.



2 - **Supervised** training on a specific task with a labeled dataset.

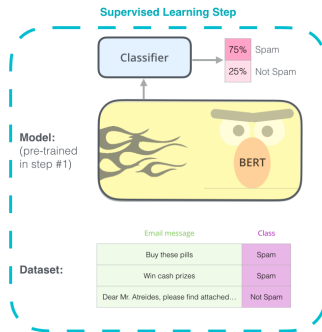


Figure: Using BERT for Transfer Learning (Alammar, 2018b)



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

The BERT model

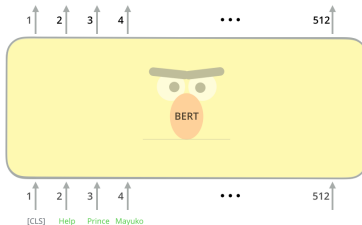


Figure: The BERT model (Alammar, 2018b)



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

BERT Architecture

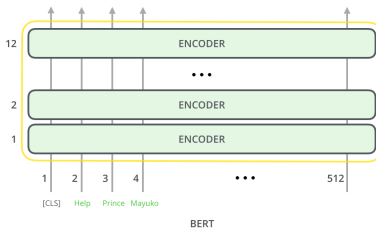


Figure: Opening up BERT (Alammar, 2018b)



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

Training Task 1: Masked Language Model

Use the output of the masked word's position to predict the masked word

Possible classes:
All English words

0.1%	Aardvark
...	...
10%	Improvisation
...	...
0%	Zyzyva

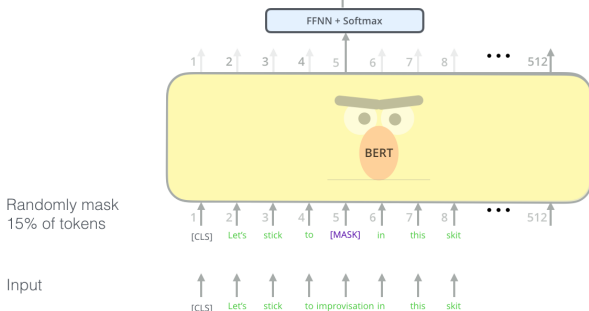


Figure: Masked Language Modeling (Alammar, 2018c)



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

Training task 2: Next Sentence Prediction

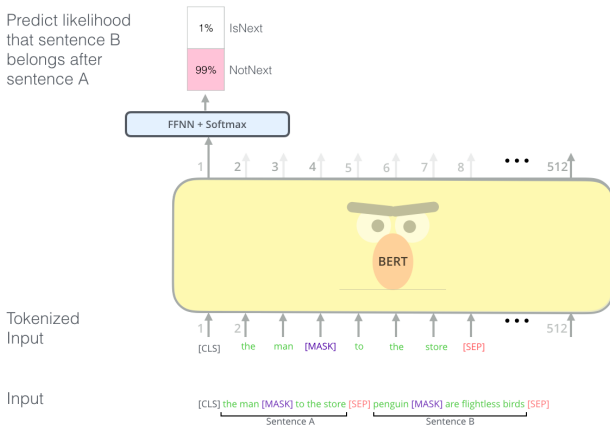


Figure: Next Sentence Prediction (Alammar, 2018c)



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

Using BERT for Classification

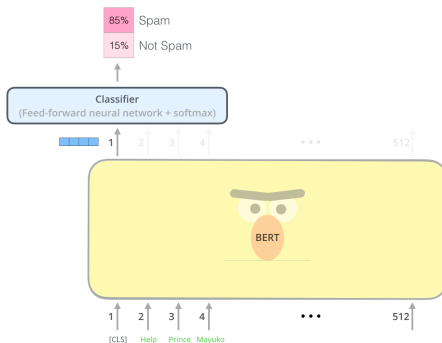


Figure: Using BERT for classification (Alammar, 2018c)



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

BERT and Contextualized embeddings

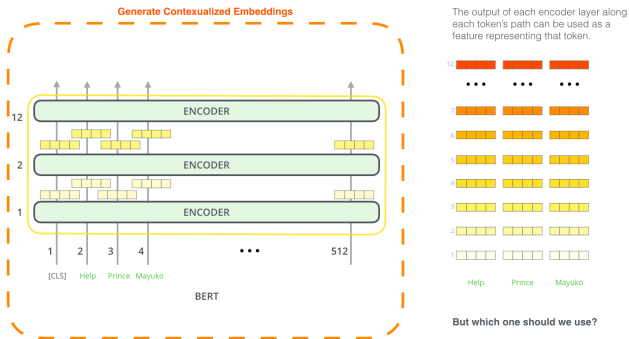


Figure: Contextualized Embeddings (Alammar, 2018c)



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

Using Contextualized Embeddings

What is the best contextualized embedding for “Help” in that context?
For named-entity recognition task CoNLL-2003 NER

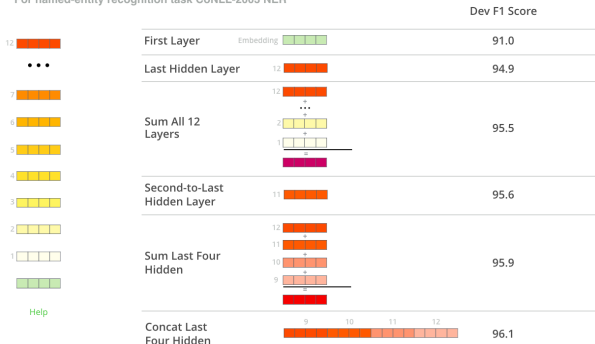


Figure: Using Contextualized Embeddings (Alammar, 2018c)



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

- RoBERTa uses a similar model as BERT with some important modifications
 1. Dynamic masking instead of Static masking



- Introduction
 - (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
 - The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
 - Transformer-Encoder Models
 - BERT
 - RoBERTa
- RoBERTa uses a similar model as BERT with some important modifications
 1. Dynamic masking instead of Static masking
 2. Longer sequences included compared to BERT



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

- RoBERTa uses a similar model as BERT with some important modifications
 1. Dynamic masking instead of Static masking
 2. Longer sequences included compared to BERT
 3. No next sentence prediction



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

- RoBERTa uses a similar model as BERT with some important modifications
 1. Dynamic masking instead of Static masking
 2. Longer sequences included compared to BERT
 3. No next sentence prediction
 4. Increased vocabulary size



- Introduction
- (Recap) Feed-Forward Neural Networks
 - Hidden Units
 - Architecture design
- The Transformer
 - Attention
 - Multi-Head Attention
 - Positional encoding
 - Add and Normalize
- Transformer-Encoder Models
 - BERT
 - RoBERTa

- RoBERTa uses a similar model as BERT with some important modifications
 1. Dynamic masking instead of Static masking
 2. Longer sequences included compared to BERT
 3. No next sentence prediction
 4. Increased vocabulary size
 5. Trained on more data (160Gb vs. 13 Gb), for longer and with larger batch sizes