

# Contextual word representations: Transformers

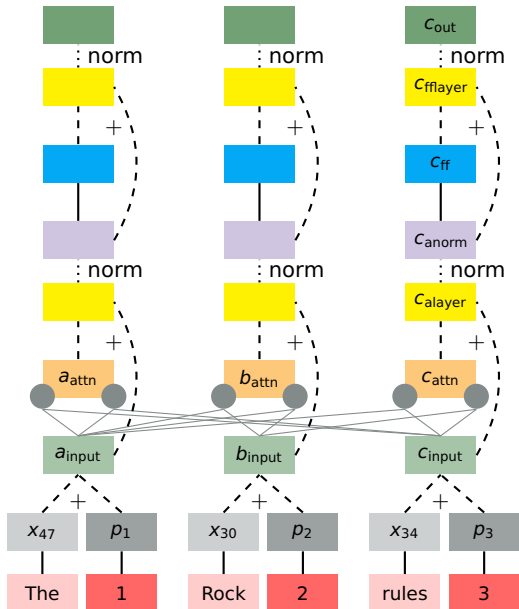
Christopher Potts

Stanford Linguistics

CS224u: Natural language understanding



# Core model structure



$$C_{out} = \frac{C_{fflayer} - \text{mean}(C_{fflayer})}{\text{std}(C_{fflayer}) + \epsilon}$$

$$C_{fflayer} = C_{anorm} + \text{Dropout}(C_{ff})$$

$$C_{ff} = \text{ReLU}(C_{anorm} W_1 + b_1) W_2 + b_2$$

$$C_{anorm} = \frac{c_{alayer} - \text{mean}(c_{alayer})}{\text{std}(c_{alayer}) + \epsilon}$$

$$C_{alayer} = \text{Dropout}(c_{attn} + c_{input})$$

$$c_{attn} = \text{sum}([\alpha_1 a_{input}, \alpha_2 b_{input}])$$

$$\alpha = \text{softmax}(\tilde{\alpha})$$

$$\tilde{\alpha} = \left[ \frac{C_{input}^T a_{input}}{\sqrt{d_k}}, \frac{C_{input}^T b_{input}}{\sqrt{d_k}} \right]$$

$$C_{input} = X_{34} + p_3$$

# Computing the attention representations

Calculation as previously given

$$\begin{aligned}c_{\text{attn}} &= \mathbf{sum}([\alpha_1 \mathbf{a}_{\text{input}}, \alpha_2 \mathbf{b}_{\text{input}}]) \\ \alpha &= \mathbf{softmax}(\tilde{\alpha}) \\ \tilde{\alpha} &= \left[ \frac{c_{\text{input}}^\top \mathbf{a}_{\text{input}}}{\sqrt{d_k}}, \frac{c_{\text{input}}^\top \mathbf{b}_{\text{input}}}{\sqrt{d_k}} \right]\end{aligned}$$

Matrix format

$$\mathbf{softmax} \left( \frac{c_{\text{input}} \begin{bmatrix} \mathbf{a}_{\text{input}} \\ \mathbf{b}_{\text{input}} \end{bmatrix}^\top}{\sqrt{d_k}} \right) \begin{bmatrix} \mathbf{a}_{\text{input}} \\ \mathbf{b}_{\text{input}} \end{bmatrix}$$

# Computing the attention representations

```
[1]: import numpy as np
```

```
[2]: seq_length = 3  
     d_k = 4
```

```
[3]: inputs = np.random.uniform(size=(seq_length, d_k))  
     inputs
```

```
[3]: array([[0.31436922, 0.66969307, 0.270804 , 0.72023504],  
          [0.87180132, 0.27637445, 0.43091867, 0.34138704],  
          [0.20292054, 0.6345131 , 0.01058343, 0.22846636]])
```

```
[4]: a_input = inputs[0]  
     b_input = inputs[1]  
     c_input = inputs[2]
```

# Computing the attention representations

```
[5]: def softmax(X):
      z = np.exp(X)
      return (z / z.sum(axis=0)).T

[6]: c_alpha = softmax([
      (c_input.dot(a_input) / np.sqrt(d_k)),
      (c_input.dot(b_input) / np.sqrt(d_k))])

[7]: c_attn = sum([c_alpha[0]*a_input, c_alpha[1]*b_input])
      c_attn

[7]: array([0.57768027, 0.48390338, 0.34643646, 0.54128076])

[8]: ab = inputs[: -1]

[9]: softmax(c_input.dot(ab.T) / np.sqrt(d_k)).dot(ab)

[9]: array([0.57768027, 0.48390338, 0.34643646, 0.54128076])

[10]: # If we allow every input to attend to itself:
      softmax(inputs.dot(inputs.T) / np.sqrt(d_k)).dot(inputs)

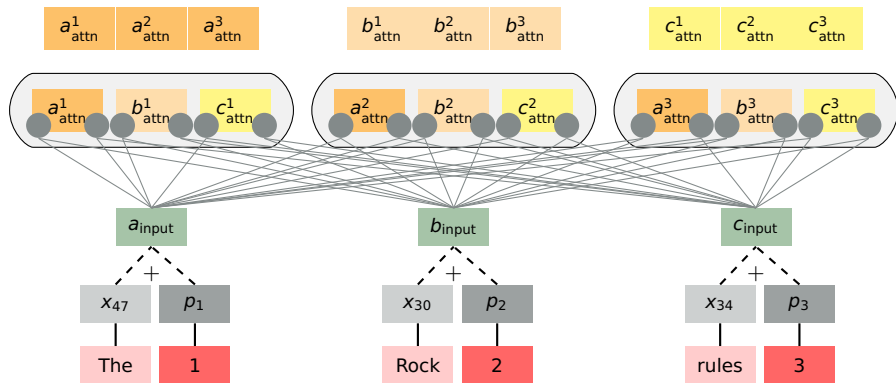
[10]: array([[0.4614388 , 0.53204444, 0.2451212 , 0.45136127],
            [0.50173123, 0.50618272, 0.26184404, 0.43678288],
            [0.45493467, 0.5332328 , 0.23643403, 0.4388242 ]])
```

# Multi-headed attention

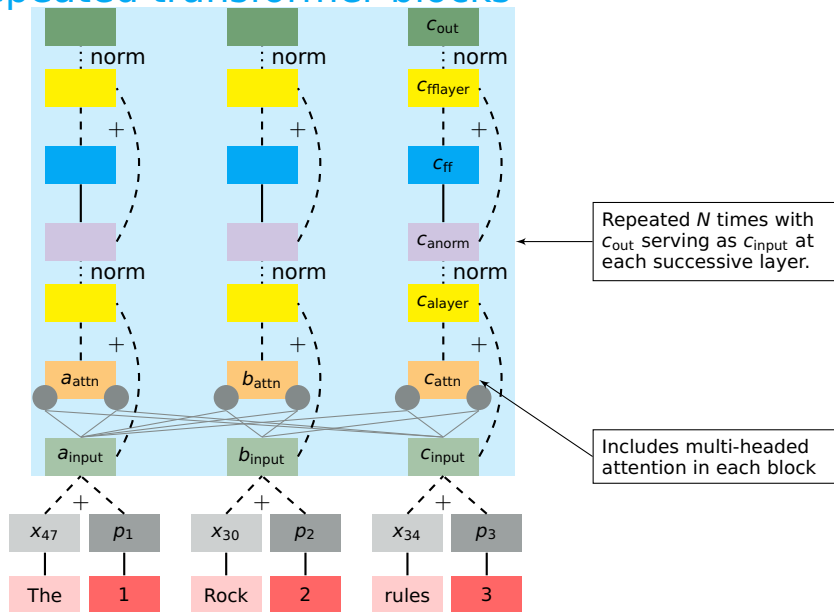
$$c_{\text{attn}}^3 = \mathbf{sum} \left( \left[ \alpha_1 (a_{\text{input}} W_3^V), \alpha_2 (b_{\text{input}} W_3^V) \right] \right)$$

$$\alpha = \mathbf{softmax}(\tilde{\alpha})$$

$$\tilde{\alpha} = \left[ \frac{(c_{\text{input}} W_3^O)^T (a_{\text{input}} W_3^K)}{\sqrt{d_k}}, \frac{(c_{\text{input}} W_3^O)^T (b_{\text{input}} W_3^K)}{\sqrt{d_k}} \right]$$



# Repeated transformer blocks



# The architecture diagram

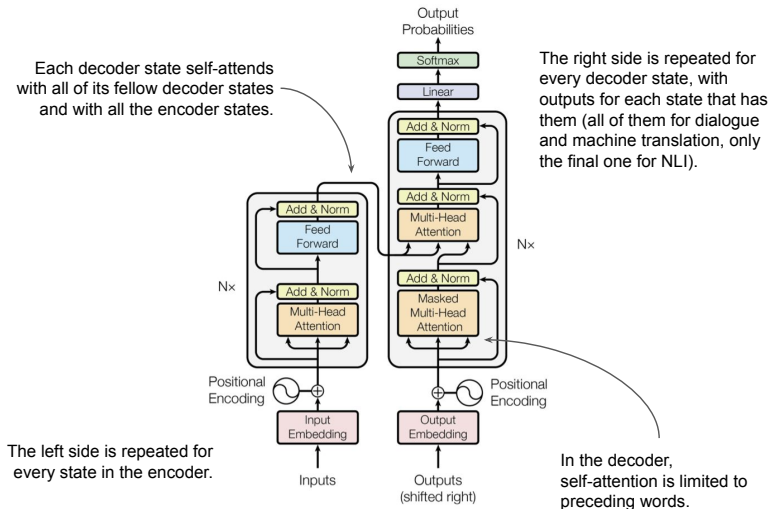


Figure 1: The Transformer - model architecture.



# References I

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. [Attention is all you need](#). In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 5998–6008. Curran Associates, Inc.