

# TTIC 31230, Fundamentals of Deep Learning

David McAllester

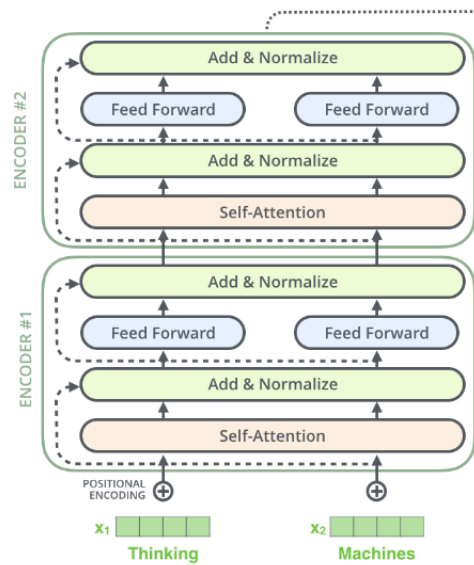
## The Transformer Part I

# The Transformer

Attention is All You Need, Vaswani et al., June 2017

The Transformer has now essentially replaced RNNs and is now used in speech, protein folding and vision.

# Vector Sequences



Each layer in the Transformer has shape  $L[T, J]$  where  $t$  ranges over the position in the input sequence and  $j$  ranges over neurons at that position (and omitting the batch index).

This is the same shape as layers in an RNN — a sequence of vectors  $L[t, J]$ .

## Parallel Layer Computation

However, in the transformer we can compute the layer  $L_{\ell+1}[T, J]$  from  $L_{\ell}[T, J]$  in parallel.

This is an important difference from RNNs which compute sequentially over time.

In this respect the transformer is more similar to a CNN than to an RNN.

# Self-Attention

The fundamental innovation of the transformer is the self-attention layer.

For each position  $t$  in the sequence we compute an attention over the other positions in the sequence.

## Transformer Heads

There is an intuitive analogy between the Transformer’s self attention and a dependency parse tree.

In a dependency parse consists of edges between words labeled with grammatical roles such as “subject-of” or “object-of”.

The self attention layers of the transformer we have “heads” which can be viewed as labels for dependency edges.

Self attention constructs a tensor  $\alpha[k, t_1, t_2]$  — the strength of the attention weight (edge weight) from  $t_1$  to  $t_2$  with head (label)  $k$ .

## Query-Key Attention

For each head  $k$  and position  $t$  we compute a key vector and a query vector with dimension  $I$  typically smaller than dimension  $J$ .

$$\text{Query}_{\ell+1}[k, t, I] = W_{\ell+1}^Q[k, I, J]L_{\ell}[t, J]$$

$$\text{Key}_{\ell+1}[k, t, I] = W_{\ell+1}^K[k, I, J]L_{\ell}[t, J]$$

$$\alpha_{\ell+1}[k, t_1, t_2] = \text{softmax}_{t_2} \frac{1}{\sqrt{I}} \text{Query}_{\ell+1}[k, t_1, I] \text{Key}_{\ell+1}[k, t_2, I]$$

## Computing the Output

$$\text{Value}_{\ell+1}[k, t, I] = W_{\ell+1}^V[k, I, J]L_{\ell}[t, J]$$

$$h_{\ell+1}^1[k, t, I] = \alpha[k, t, T]\text{Value}[k, T, I]$$

$$h_{\ell+1}^2[t, C] = h_{\ell+1}^1[0, t, I]; \cdots ; h_{\ell+1}^1[K - 1, t, I]$$

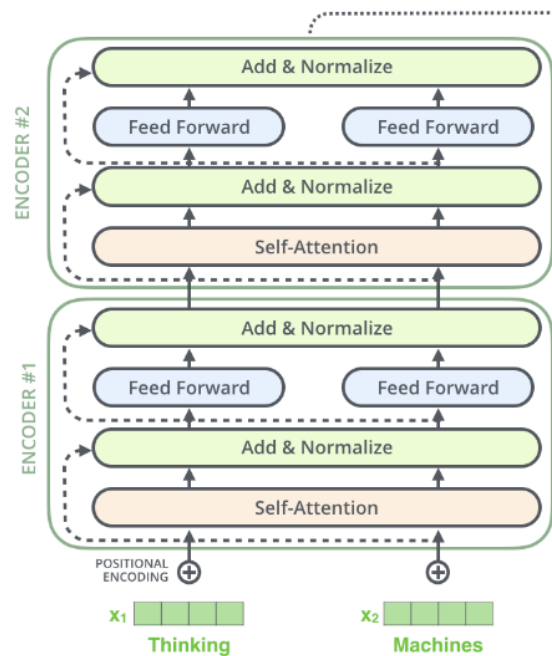
$$L_{\ell+1}[t, J] = W_{\ell+1}^0[J, C]h^2[t, C]$$

Here semicolon denotes vector concatenation.



# The Transformer Layer

Each “transformer layer” consists of six “sublayers” the first of which is the self-attention layer.



Jay Alammar's blog

The other layers are discussed in the next unit.

**END**