# Seq2Seq

## Encoder RNN

summary of what encoder has seen so far

$u_1$ $u_2$ $(h_3)$

Output condensed representation

$x_1$ $x_2$ $x_2$

each vector $x$ is embedding of a word
English sentence

## Decoder

$S_1$ $S_2$ $S_3$

$S_0$

$x_1'$ $x_2'$ $x_3'$

generated word

too many words will lead to forgetting in $(h_t)$ of encoder

using attention

BLEU

20 words    50 words    seq length

translation quality

## Attention for Seq2Seq

- encoder won't forget
- $O(N^2)$ time complexity

- applied as decoder starts working

From encoder: $u_1$ $h_2$ $h_3$ ... $h_m = S_0$ in decoder

weight $= \alpha_i = align(h_i, S_0)$ how well $h_i$ and $S_0$ match 0-1

↑ weight $=$ ↑ relevant

$\sum_m \alpha_i = 1$

matrix, trainable, shared

$$\widetilde{\alpha_i} = v^T \cdot tanh\left[ W \cdot \binom{h_i}{S_0} \right]$$

vector, trainable, shared

normalize: $\alpha_i = softmax(\widetilde{\alpha_i})$ so add to 1

— Another align function

$$k_i = W_k \cdot h_i \quad \text{for } i = 1 \dots m$$

$$q_0 = W_Q \cdot s_0$$

$$\overline{\alpha_i} = k_i^T \cdot q_0 \quad \text{for } i = 1 \dots m$$

$$\alpha_i = \text{softmax}(\overline{\alpha_i})$$

— Context vector

$$c_0 = \alpha_1 \cdot h_1 + \dots + \alpha_m \cdot h_m \quad \text{\color{blue}weighted avg.}$$

$\color{blue}\hookrightarrow$ corresponds to $s_0$

<span style="color:red">Simple RNN:</span>
$$s_1 = \tanh\left(A^1 \cdot \begin{pmatrix} x_1^1 \\ s_0 \end{pmatrix} + b\right)$$

$$x_1^1 = \text{<START>} \text{ in machine translation}$$

<span style="color:red">Simple RNN + Attention:</span>
$$s_1 = \tanh\left(A^1 \cdot \begin{pmatrix} x_1^1 \\ s_0 \\ c_0 \end{pmatrix} + b\right)$$
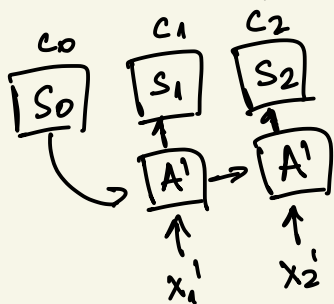
• compute $c_1$ corresponding to $s_1$

$$\alpha_i = \text{align}(h_i, s_1) \text{ weights} \quad \color{blue}\text{how well } h_i \text{ and } s_1 \text{ are aligned (all encoder hidden states + decoded vector)}$$

• $\alpha$ is recalculated on each run
  for $s_0 \dots s_n$
$$c_1 = \alpha_1 h_1 + \dots + \alpha_m \cdot h_m$$

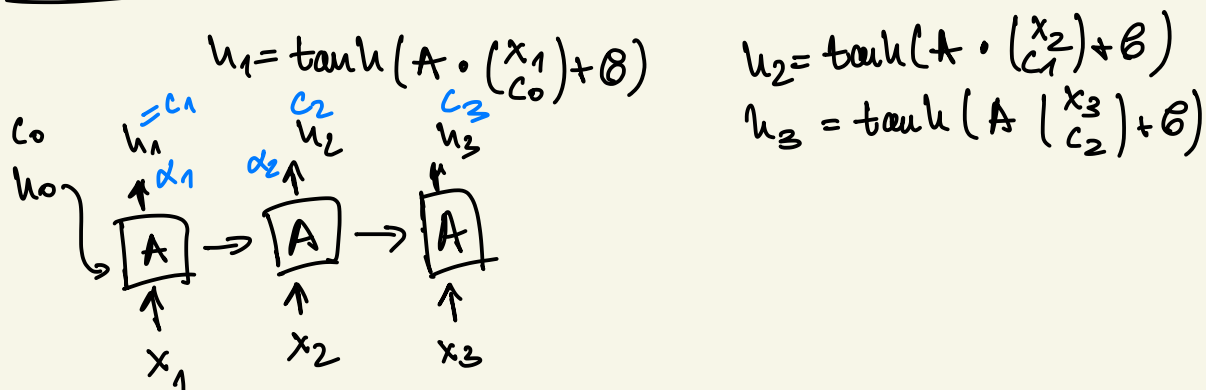• $s_2 = \tanh\left(A^1 \cdot \begin{pmatrix} x_2^1 \\ s_1 \\ c_1 \end{pmatrix} + b\right)$



• Complexity:  $c_j$ computed from $m$ weights $\alpha_1 \dots \alpha_m$
  w.t weights      $m =$ encoder steps
                   $t =$ decoder steps

without attention: $O(m + t)$
with: $O(m \cdot t)$

# Self Attention for RNN

$$h_1 = \tanh\left(A \cdot \begin{pmatrix} x_1 \\ c_0 \end{pmatrix} + \theta\right) \qquad h_2 = \tanh\left(A \cdot \begin{pmatrix} x_2 \\ c_1 \end{pmatrix} + \theta\right)$$

$$h_3 = \tanh\left(A \begin{pmatrix} x_3 \\ c_2 \end{pmatrix} + \theta\right)$$



Weights: $\alpha_i = \text{align}(h_i, h_2)$   $i = $ each existing state

Ex.: $c_3 = \alpha_1 h_1 + \alpha_2 h_2 + \alpha_3 h_3$

- With self attention RNN is less likely to forget
- Self attention is applied to a simple RNN (not Seq2Seq model)
- Pay attention to relevant context in the past

# Transformer Model (Attention Layer)

- 2 years after Seq2Seq attention
- 1 year after self attention
- Seq2Seq, Encoder-Decoder, pure dense layers (no RNN)
- higher accuracy than RNN

- Encoder weights   $\alpha_{ij} = \text{align}(h_i, s_j)$   $j = $ decoder hidden state at step $j$

  for $i = 1..m$   $i = $ encoder state at step $i$

— Calculate align:
- $k_i = W_k \cdot h_i$   $W_k = $ param matrix, random, learns from training data
- $k = [k_1 \cdots k_m]$
- $q_{:j} = W_Q \cdot s_j$
- $\alpha_{:j} = \text{softmax}(k^T q_{:j})$   $k^T$ state vec transformed
- $\alpha_{:j} = \begin{pmatrix} \alpha_{1j} \\ \vdots \\ \alpha_{mj} \end{pmatrix} = \text{softmax}\left(\begin{bmatrix} = \\ = \\ = \end{bmatrix} \begin{bmatrix} \vdots \\ q_{:j} \end{bmatrix}\right)$

- In transformer:
  - $q_j$ = Query (to match others) = $W_Q \cdot s_j$
  - $k_i$ = key (to be matched) = $W_k \cdot h_i$   W to be learnt
  - $v_i$ = Value (to be weighted avg) = $W_v \cdot h_i$

  Context vector: $c_j = \alpha_{1j} \cdot v_1 + \dots + \alpha_{mj} \cdot v_m$

- Intuition

  $s_j$ = decoder output at step $j$

  - Query : $q_j = W_Q \cdot s_j$
  - key : $k_i = W_k \cdot h_i$ ~ for all encoder hidden states: $= k$
  - Weight : $\alpha_j = $ softmax $(k^T q_j)$
  - Value : $v_i = W_v \cdot h_i$ for all encoder states (m weights)
  - Context vect : $c_j = \alpha_{1j} v_1 + \dots + \alpha_{mj} v_m$

— Attention without RNN
  - Encoder : $x_1 \dots x_m$ english words
  - Decoder : $x'_1 \dots x'_t$ german words
  - start decoder with <start>, continue until receiving <stop>
  - key = $k_i = W_K \cdot x_i$
  - Value = $v_i = W_v \cdot x_i$

  $\underset{X_1}{\overset{k_1 \ v_1}{\uparrow}}$   $\underset{X_2}{\overset{k_2 \ v_2}{\uparrow}}$   $\underset{X_3}{\overset{k_3 \ v_3}{\uparrow}}$ .... Encoder side

  - Query = $q_j = W_Q \cdot x'_j$

  $\underset{X'_1}{\overset{q_1}{\uparrow}}$   $\underset{X'_2}{\overset{q_2}{\uparrow}}$   $\underset{X'_3}{\overset{q_3}{\uparrow}}$ .... Decoder side

  - Parameter matrices $W_K / W_v / W_Q$
  - Compute weights: $\alpha_1 = $ Softmax $(k^T \cdot q_1)$
  - Compute context: $c_1 = \alpha_{11} \cdot v_1 + \dots \alpha_{m1} v_m = V \alpha_1$

- $c_j$ is function of $x'_j$ and $x_1 \cdots x_m$
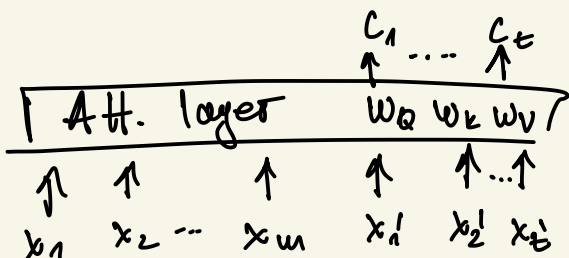
$c_2 \Rightarrow$ softmax classifier $\Rightarrow P_2$ <span style="color:red">$\rightarrow$ sample German word $\rightarrow x'_3 =$ Embedding of selected word</span>
<span style="color:blue">$\hookrightarrow$ knows all english</span> <span style="color:blue">$\hookrightarrow$ probability distribution over German Vocab</span>
<span style="color:blue">+ previous german words</span>
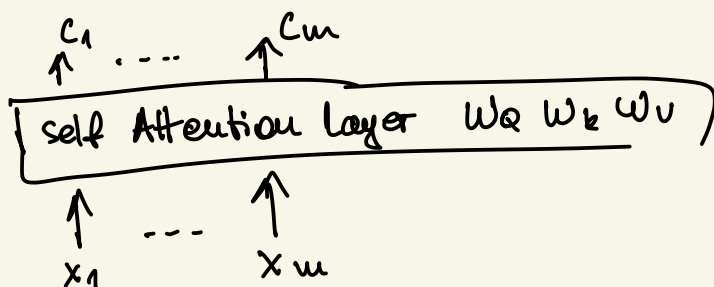
- Attention Layer: $[c_1 \cdots c_t] = Attn(x, x')$

  Encoder inputs: $x = x_1 \cdots x_m$
  Decoder inputs: $x' = x'_1 \cdots x'_t$
  Parameters: $W_Q \ W_K \ W_V$

$$c_1 \cdots c_t$$
$$\uparrow \qquad \uparrow$$

| Att. layer $\quad W_Q \ W_K \ W_V$ |
$$\uparrow \ \uparrow \qquad \uparrow \qquad \uparrow \quad \uparrow \cdots \uparrow$$
$$x_1 \ x_2 \cdots x_m \qquad x'_1 \ x'_2 \ x'_t$$

— Self Attention

$$c_1 \cdots \qquad c_m$$
$$\uparrow \qquad \qquad \uparrow$$

{ self Attention Layer $\quad W_Q \ W_K \ W_V$ )
$$\uparrow \quad \cdots \quad \uparrow$$
$$x_1 \qquad \quad x_m$$

$$q_i = W_Q \cdot x_i \qquad k_i = W_K \cdot x_i \qquad v_i = W_V x_i$$

$q_1 \ k_1 \ v_1 \qquad q_2 \ k_2 \ v_2 \qquad q_3 \ \frac{k_3}{\uparrow} \ v_3 \qquad \cdots \ \frac{q_m \ k_m \ v_m}{\uparrow}$
$$\uparrow \qquad\qquad \uparrow \qquad\quad x_3 \qquad\qquad x_m$$
$$x_1 \qquad\qquad x_2$$

Weights: $\alpha_{;j} = softmax(k^T \cdot q_j) \qquad j = 1 \cdots m$

Context vect: $c_1 = \alpha_{11} \cdot v_1 + \cdots + \alpha_{m1} v_m = V \cdot \alpha_1$
$$c_{;j} = \alpha_{1j} v_1 + \cdots + \alpha_{mj} v_m = V \alpha_{;j}$$
$$\Rightarrow c_{;j} = V \cdot softmax(k^T \cdot q_{;j})$$

Self attention $[c_1 \cdots c_m] = Attn(x, x)$ $\qquad$ Inputs: $x_1 \cdots x_m$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ Params: $W_Q \ W_K \ W_V$
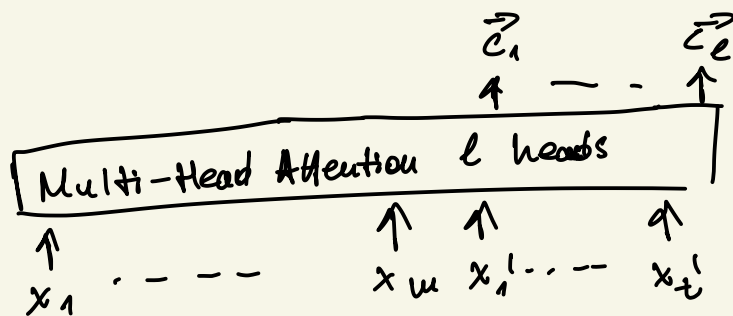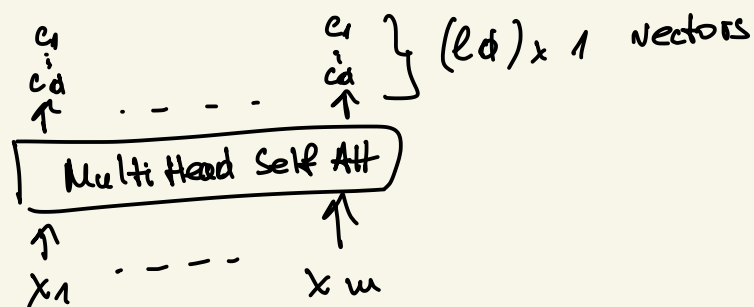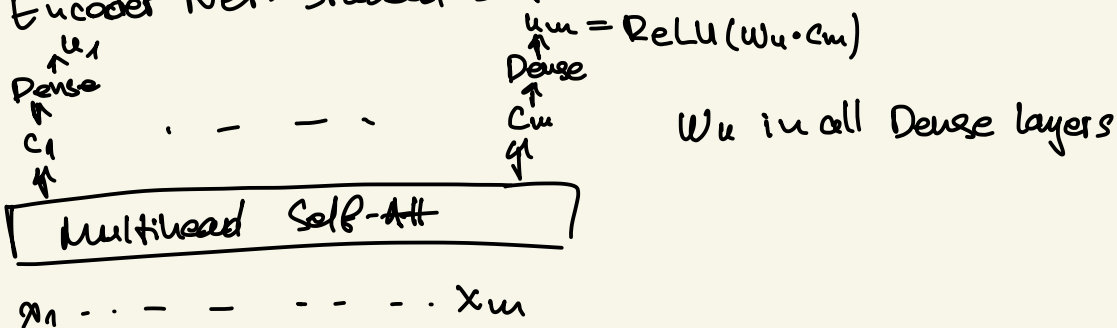
- $c_j$ is influenced by all $x_1 \ldots x_m$
- Attention for Seq2Seq RNN 2015
  Neural machine translation by jointly learning to align and translate
- Self attention for all RNN (not only Seq2Seq) 2016
  LSTM Networks for machine reading
- Attention without RNN 2017
  Attention is all you need
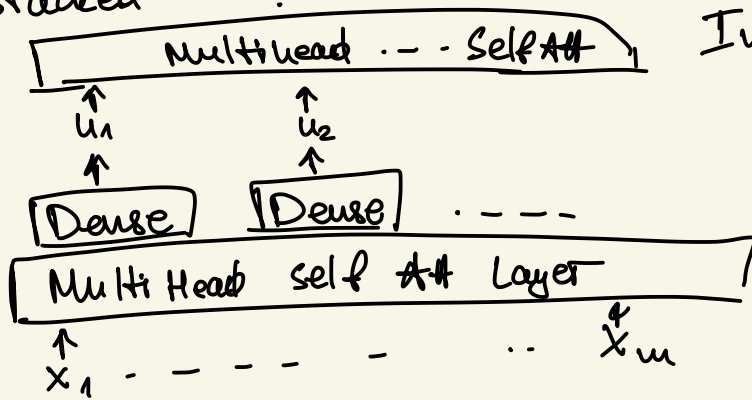
## Transformer Model: Deep Neural Network

- Attention / self attention = single head attention

— Multihead Self-Attention
  - single head params: $W_Q$ $W_k$ $W_v$
  - $l$ single head attention (not sharing matrices) has $3l$ matrices
  - concatenate context outputs of single head
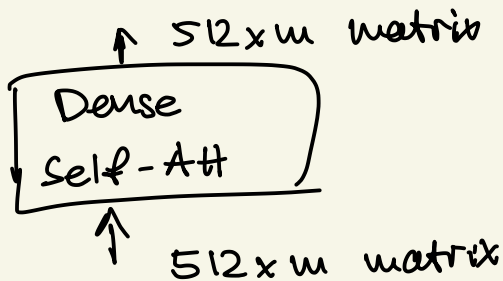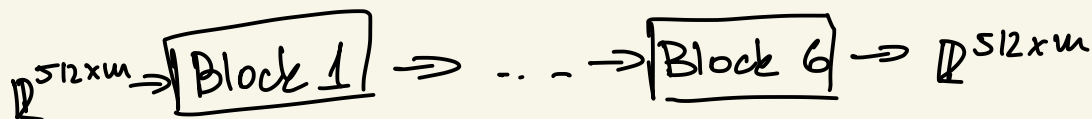


— Encoder Net: Stacked Self-Attention



$u_m = ReLU(W_u \cdot c_m)$

$W_u$ in all Dense layers

— Stacked



$Input = Output \; len = m$

— Block of Encoder



512 can be different
(embedding dependant)

$\mathbb{R}^{512 \times m} \Rightarrow \boxed{Block \; 1} \Rightarrow \cdots \Rightarrow \boxed{Block \; 6} \Rightarrow \mathbb{R}^{512 \times m}$
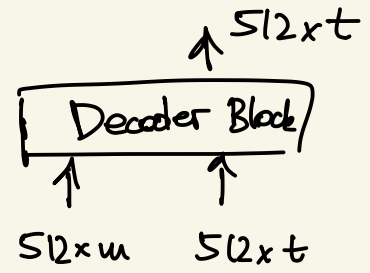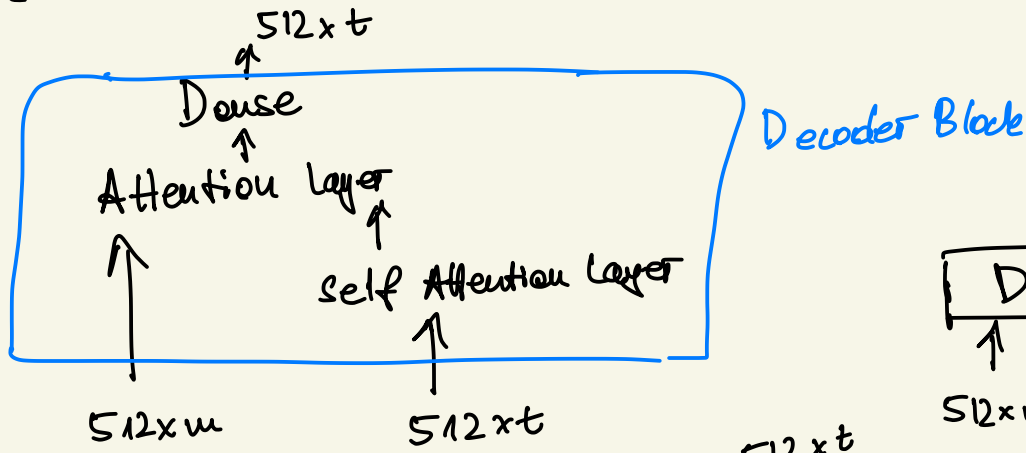
— Decoder Net

- Transformer is seq2seq (encoder + decoder) model
- Encoder input: embeddings $x_1 \cdots x_m$
- Decoder inputs: $x_1' \cdots x_t'$ generated $t$ words
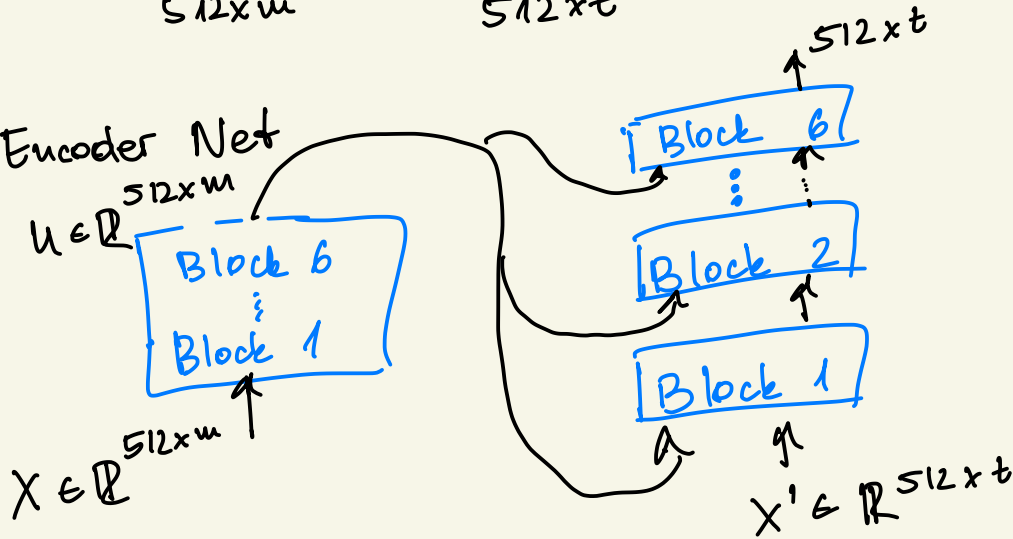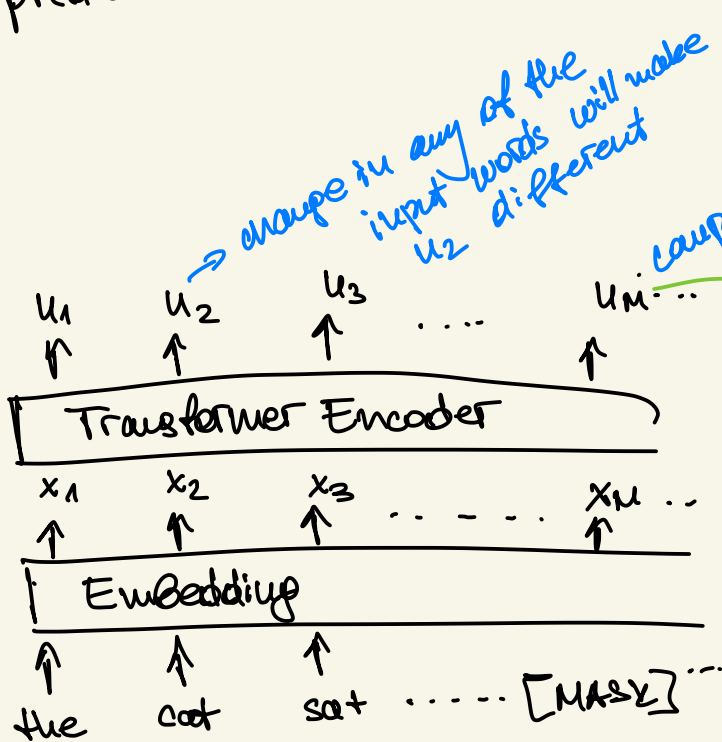- 1 block $\approx$ 1 multi-head attention layer + 1 dense



all inputs 512 length

- ## Decoder Block

$512 \times t$

Dense

Attention Layer

Self Attention Layer

Decoder Block

$512 \times m$          $512 \times t$

$512 \times t$

Decoder Block

$512 \times m$     $512 \times t$

- Encoder Net

$u \in \mathbb{R}^{512 \times m}$

| Block 6 |
| Block 1 |

$X \in \mathbb{R}^{512 \times m}$

$512 \times t$

Block 6
Block 2
Block 1

$X' \in \mathbb{R}^{512 \times t}$

## Bert for pretrained transformers

- predict masked word
- predict next sentence

→ change in any of the input words will make $u_2$ different

$u_i$ computed with u context

$u_1$  $u_2$  $u_3$  ....  $u_M$ ...

→ softmax classifier → $p$ distribution
Loss = CrossEntropy $(e, p)$
gradient descent to update params

Transformer Encoder

$x_1$  $x_2$  $x_3$  - - - - -  $x_M$ ..

Embedding

the    cat    sat   ..... [MASK] ...

— Predict next sentence
- Input [CLS] sentence 1 [SEP] sentence 2
  ↳ separation
  ↳ for classification

$c$ $u_1 \ldots u_m$ $S$ $u_1 \ldots u_t$     $c = 0$   not real sentence 2

Encoder
Embedding
[CLS]   sentence 1 [SEP]   sentence 2

- improves embedding layer
- improves attention on how well words are aligned
- labels are automatically generated

## Vision Transformer

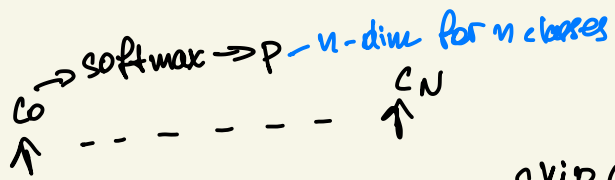image $\Rightarrow$ $\boxed{NN}$ $\rightarrow$ $\vec{P}$   $n$-dimensional for $n$ classes (confidences)

- vision transformer beats CNN 2021

- partition image into patches of same shape
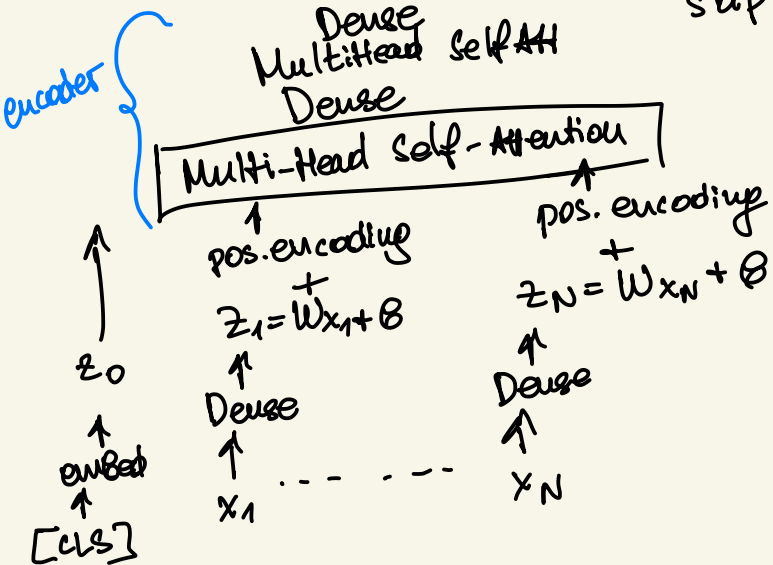  patch size = 16×16 default
  stride = non overlapping patches 16×16 default

- reshape tensors (3 values for RGB) into vectors    $d_1 \times d_2 \times d_3$ tensor
  ↳ $x_1 \ldots x_n$ $\Rightarrow$ $d_1 \cdot d_2 \cdot d_3 \times 1$ vector

↳ softmax $\Rightarrow$ P — $n$-dim for $n$ classes
$c_0$ ----------- $c_N$

skip connection + normalization

encoder {
Dense
MultiHead Self-Att
Dense
$\boxed{\text{Multi-Head Self-Attention}}$

pos.encoding                 pos.encoding
+                            +
$z_1 = W x_1 + \theta$        $z_N = W x_N + \theta$   $\theta / W$ are shared

$z_0$         Dense              Dense
              $x_1$ --- --- $x_N$
embed
[CLS]

– Training

Random Initialize $\rightarrow$ pretrained $\Rightarrow$ Fine Tuned $\longrightarrow$ Test accuracy

Dataset A        Dataset B        Test set B

very large       smaller for target action

– Datasets

Imagenet small   1,3M      ResNet     ViT

Imagenet med   14M        |———|———————|→

JFT        300M           100M    300M