

# GESTURE GENERATION USING TRANSFORMER DECODER ARCHITECTURE

**Motive :** Given (Having predicted)  $t$  clips each containing  $T$  frames along with the word embeddings,  $\Phi_t$  from a pre-trained model, predicting  $(t+1)^{\text{th}}$  clip, with the help of self-attention and attention w.r.t word embeddings.

## Attention Block

The attention block consists of -

1. Masked Separable self-attention block
2. Frame level attention block with text embeddings

## Masked Separable self-attention block

Let the real video be represented as  $v$ , with dimensions  $H \times W \times N \times C$ , where  $H$  is the height and  $W$  is the width and  $C$  is the number of channels, and  $N$  is the total number of frames in the video clip. The whole video is divided into  $K$  different clips of  $T$  frames each, where  $v \in \mathbb{R}^{H \times W \times N \times C}$ , and  $v^t \in \mathbb{R}^{H \times W \times T \times C}$

Separable attention operation is performed across Time, Height and Width

### Across Time

Video clip  $v^j$  is reshaped to  $v_n^j$  such that  $v_n^j \in \mathbb{R}^{(H \times W) \times T \times C}$ , and with the current clip,  $v_n^t$  as query, and the other video clips being keys and values, it is attended with the other masked  $(K-t+1)$  blocks and unmasked  $t$  blocks of the video frames.

For each  $j \in \{0, 1, \dots, K-1\}$  and for each  $i \in \{0, 1, \dots, T-1\}$ ,

$$Q_1^i = v_n^t(i).q_1 \quad \dots (1)$$

$$K_1^i = v_n^j(i).k_1 \quad \dots (2)$$

$$V_1^i = v_n^j(i).v_1 \quad \dots (3)$$

Where,  $v_n^j(i)$  is the  $i^{\text{th}}$  2D matrix of  $v_n^j$  across time dimension, and  $\{q_1, k_1, v_1\} \in \mathbb{R}^{C \times C'}$  are shared weights for all the  $T$  frames and  $\{Q_1^i, K_1^i, V_1^i\} \in \mathbb{R}^{H \times W \times C'}$

$$A_1^i = \text{softmax}(Q_1^i \cdot (K_1^i)^T) \cdot V_1^i \quad \dots (4)$$

$$A_1^i = A_1^i \cdot (W_1^0)^T \quad \dots (5)$$

Where,  $W_1^0 \in \mathbb{R}^{C \times C}$

Finally,

the attention output is given as  $A_1 = [A_1^0, A_1^1, \dots, A_1^{T-1}]$ , such that  $A_1 \in \mathbb{R}^{H \times W \times T \times C}$

The Video frames are hence updated as  $v_n^t = A_1$ . This is then reshaped and used for attention across height, followed by width.

## Across Height

Video clip  $v_n^j$  is reshaped to  $v_h^j$  such that  $v_h^j \in \mathbb{R}^{(W \times T) \times H \times C}$ , and with the current clip,  $v_h^t$  as query, and the other video clips being keys and values, it is attended with the other masked  $(K-t+1)$  blocks and unmasked  $t$  blocks of the video frames.

For each  $j \in \{0, 1, \dots, K-1\}$  and for each  $i \in \{0, 1, \dots, H-1\}$ ,

$$Q_2^i = v_h^t(i).q_2 \quad \dots (6)$$

$$K_2^i = v_h^j(i).k_2 \quad \dots (7)$$

$$V_2^i = v_h^j(i).v_2 \quad \dots (8)$$

Where,  $v_h^j(i)$  is the  $i^{\text{th}}$  2D matrix of  $v_h^j$  across height dimension, and  $\{q_2, k_2, v_2\} \in \mathbb{R}^{C \times C}$  are shared weights for all the  $T$  frames and  $\{Q_2^i, K_2^i, V_2^i\} \in \mathbb{R}^{W \times T \times C}$

$$A_2^i = \text{softmax}(Q_2^i \cdot (K_2^i)^T) \cdot V_2^i \quad \dots (9)$$

$$A_2^i = A_2^i \cdot (W_2^0)^T \quad \dots (10)$$

Where,  $W_2^0 \in \mathbb{R}^{C \times C}$

Finally,

the attention output is given as  $A_2 = [A_2^0, A_2^1, \dots, A_2^{H-1}]$ , such that  $A_2 \in \mathbb{R}^{W \times T \times H \times C}$

The Video frames are hence updated as  $v_h^t = A_2$ . This is then reshaped and used for attention across width.

## Across Width

Video clip  $v_n^j$  is reshaped to  $v_w^j$  such that  $v_w^j \in \mathbb{R}^{(H \times T) \times W \times C}$ , and with the current clip,  $v_w^t$  as query, and the other video clips being keys and values, it is attended with the other masked  $(K-t+1)$  blocks and unmasked  $t$  blocks of the video frames.

For each  $j \in \{0, 1, \dots, K-1\}$  and for each  $i \in \{0, 1, \dots, W-1\}$ ,

$$Q_3^i = v_w^t(i) \cdot q_3 \quad \dots (11)$$

$$K_3^i = v_w^j(i) \cdot k_3 \quad \dots (12)$$

$$V_3^i = v_w^j(i) \cdot v_3 \quad \dots (13)$$

Where,  $v_w^j(i)$  is the  $i^{\text{th}}$  2D matrix of  $v_w^j$  across width dimension, and  $\{q_3, k_3, v_3\} \in \mathbb{R}^{C \times C'}$  are shared weights for all the  $T$  frames and  $\{Q_3^i, K_3^i, V_3^i\} \in \mathbb{R}^{H \times T \times C'}$

$$A_3^i = \text{softmax}(Q_3^i \cdot (K_3^i)^T) \cdot V_3^i \quad \dots (14)$$

$$A_3^i = A_3^i \cdot (W_3^0)^T \quad \dots (15)$$

Where,  $W_3^0 \in \mathbb{R}^{C \times C'}$

Finally,

the attention output is given as  $A_3 = [A_3^0, A_3^1, \dots, A_3^{W-1}]$ , such that  $A_3 \in \mathbb{R}^{H \times T \times W \times C}$

The Video frames are hence updated as  $v_w^t = A_3$ . This is then reshaped and used for attention across width.

After all the above steps  $v_w^t$  is reshaped to  $v^{\text{temp}} \in \mathbb{R}^{H \times W \times T \times C}$ , followed by addition with  $v^t$  and Layer Normalisation as follows,

$$\begin{aligned} v_t &= (v_{\text{temp}} + v_t) \\ \mu_t &= \frac{1}{T} \sum_{i=1}^T v_t(i) \quad , \quad \sigma_t^2 = \frac{1}{T} \sum_{i=1}^T (v_t(i) - \mu_t)^2 \\ v_t(i) &= (v_t(i) - \mu_t) / \sqrt{(\sigma_t^2)} \end{aligned}$$

## Frame level attention block with text embeddings

Masked self attention is followed by this block, where each frame of  $v^t$  is attended w.r.t the word embeddings, where every frame acts a query, and word embeddings act as keys and values. Inspired from AttnGan paper.

Let word embeddings be obtained from a pretrained model and be represented as  $\Phi_t$ , such that  $\Phi_t \in \mathbb{R}^{D \times L}$ .  $\Phi_t$  is brought to the same semantic space as  $v^t$  using a perceptron layer such that,

$$e' = U \cdot \Phi_t, \text{ where } U \in \mathbb{R}^{H \times W \times D}$$

$L$  is the number of words and  $D$  is the dimensionality of the word feature vector, and  $e' \in \mathbb{R}^{H \times W \times L}$

The video frames are divided into T matrices  $v_i^t \in \mathbb{R}^{H \times W \times C}$  for  $i \in \{0, 1, \dots, T-1\}$

Let  $v_i^t(j)$  denote the  $j^{\text{th}}$  feature of  $v_i^t$ , such that  $v_i^t(j) \in \mathbb{R}^{H \times W \times 1}$  and  $j \in \{0, 1, \dots, C-1\}$

Also, let  $e'(k)$  denote the  $k^{\text{th}}$  column of  $e'$ , such that  $e'(k) \in \mathbb{R}^{H \times W \times 1}$  and  $k \in \{0, 1, \dots, D-1\}$

$$a_i(j, k) = \sum_k \text{softmax}((v_i^t(j))^T \cdot e'(k)) \cdot e'(k), \text{ such that } a_i(j, k) \in \mathbb{R}^{H \times W \times 1}$$

When all j features of frame with all the k word features, we get  $a_i \in \mathbb{R}^{H \times W \times C}$  and after all the attentions across all the frames are calculated we get,

$$a = [a^0, a^1, \dots, a^{T-1}], \text{ such that } a \in \mathbb{R}^{H \times W \times C \times T}$$

Attention output, a, is then reshaped to  $H \times W \times C \times T$ , and then added to the input  $v^t$  and then normalized, just like the previous layer

This completes the attention block.

## Generator

Each attention block will be followed by a series/(single) of convolution layer(s) to downscale the image, to bring it down to  $v_m^t \in \mathbb{R}^{h \times w \times t \times c}$ . Randomness (z) is introduced (experimental) such that,  $z \sim \mathcal{N}(0, I)$ , derived from Gaussian Distribution, such that  $z \in \mathbb{R}^{h \times w \times t \times c'}$  and concatenated with  $v_m^t$  along the channel dimension.

The generator, G is then defined as,

$$G : \{\mathbb{R}^{h \times w \times t \times c}, \mathbb{R}^{h \times w \times t \times c'}\} \rightarrow \mathbb{R}^{H \times W \times T \times C}, \text{ i.e.,}$$

$$G(z \parallel v_m^t) = v^{\text{pred}}, \text{ such that } v^G = \{f_1^G, f_2^G, \dots, f_T^G\}$$

Where,  $f_t^G \in \mathbb{R}^{H \times W \times C}$

## Discriminator

Inspired by the paper, “To create what you tell”, there are 3 different discriminator networks namely,

1. Video Discriminator ( $D_0$ )
2. Frame Discriminator ( $D_1$ )
3. Motion Discriminator ( $D_2$ )

Let s be the sentence embedding given by the pre-trained text encoder model

$$D_0(v, s) : \{\mathbb{R}^{\text{dv}}, \mathbb{R}^{\text{ds}}\} \rightarrow \{0, 1\}$$

$$D_1(f_i, s) : \{R^{df}, R^{ds}\} \rightarrow \{0, 1\}$$

$$D_2(f_i, f_{i-1}) : \{R^{df}, R^{ds}\} \rightarrow R^{c0 \times h0 \times w0}$$

Where,  $dv$  is the dimensionality of input video (either generated or ground truth),  $ds$  is the dimensionality of the sentence embedding,  $f_i$  is the  $i^{th}$  frame of the video clip and  $c0, h0, w0$  are the dimensions of the downsampled frame  $m_f^i$

## Losses

### Video-level Matching aware loss

Let  $v^+$  be the real video with the correct sentence,  $v^-$  be the real video with mismatched sentence and  $v^{pred}$  be the generated video from G. This loss is calculated across the whole video produced, i.e  $v \in R^{H \times W \times T \times C}$

$$L_v = -\frac{1}{3} [\log(D_0(v^+, s)) + \log(1 - D_0(v^-, s)) + \log(1 - D_0(v^{pred}, s))]$$

### Frame-level Matching aware loss

Let  $f^+(i)$  be the  $i^{th}$  frame from the real video with the correct sentence,  $f^-(i)$  be the  $i^{th}$  frame from the real video with mismatched sentence and  $f^{pred}(i)$  be the  $i^{th}$  frame from the generated video from G. This loss is calculated from a single frame produced, i.e  $f(i) \in R^{H \times W \times C}$

$$L_f = -\frac{1}{3N} \left[ \sum_{i=1}^N \log(D_1(f^+(i), s)) + \sum_{i=1}^N \log(1 - D_1(f^-(i), s)) + \sum_{i=1}^N \log(1 - D_1(f^{pred}(i), s)) \right]$$

### Temporal Coherence loss

Let  $m^f(i)$  be the downsampled  $i^{th}$  frame from the generated video from G. This loss is calculated from 2 consecutive frames produced, i.e  $m^f(i) \in R^{h0 \times w0 \times c0}$ . Simply calculates the Euclidean Distance between 2 consecutive frames. This loss function corresponds only to the generator.

$$L_t = \frac{1}{N-1} \sum_{i=2}^N \|m^f(i) - m^f(i-1)\|_2^2$$

**Discriminator Loss :**  $L_D = \frac{1}{2} (L_v + L_f)$

**Generator Loss :**  $L_G = -\frac{1}{3} [\log(D_0(v^{pred}, s)) + \frac{1}{N} \sum_{i=1}^N \log(D_1(f^{pred}(i), s)) - L_t]$