# TransPose: Keypoint Localization via Transformer — Supplementary Material

#### A. 2D Sine Position Embedding

Without the position information embedded in the input sequence, the Transformer Encoder is a permutationequivariant architecture:

Encoder 
$$(\rho(\mathbf{X})) = \rho(\text{Encoder}(\mathbf{X})),$$
 (1)

where  $\rho$  is any permutation for the pixel locations or the order of sequence. To make the order of sequence or the spatial structure of the image pixels matter, we follow the sine positional encodings but further hypothesize that the position information is independent at x (horizontal) and y (vertical) direction of an image, like the ways of [4, 2]. Concretely, we keep the original 2D-structure respectively with d/2 channels for x, y-direction:

$$PE_{(2i,p_y,:)} = \sin\left(2\pi * p_y/(H * 10000^{2i/\frac{d}{2}})\right),$$

$$PE_{(2i+1,p_y,:)} = \cos\left(2\pi * p_y/(H * 10000^{2i/\frac{d}{2}})\right),$$

$$PE_{(2i,:,p_x)} = \sin\left(2\pi * p_x/(W * 10000^{2i/\frac{d}{2}})\right),$$

$$PE_{(2i+1,:,p_x)} = \cos\left(2\pi * p_x/(W * 10000^{2i/\frac{d}{2}})\right),$$
(2)

where i = 0, 1, ..., d/2 - 1,  $p_x$  or  $p_y$  is the position index along x or y-direction. Then they are stacked and flattened into a shape  $\mathbb{R}^{L \times d}$ . The position embedding is injected into the input sequences before self-attention computation. We use 2D sine position embedding by default for all models.

# B. What position information has been learned in the TransPose model with learnable position embedding?

We show what position information has been learned in the TransPose (TransPose-R) with learnable position embedding. It has been discussed in the paper. As shown in Fig. 1, we visualize the similarities by calculating the cosine similarity between vectors at any pair of locations of the learnable position embedding and reshaping it into a 2D grid-like map. We find that the embedding in each location of learnable position embedding has a unique vector value in the *d*-dim vector space, but it has relatively higher cosine



Figure 1. The cosine similarities between the learned position embedding vectors, which have been reshaped into 2D grid and interpolated with 0.25 scale factor for a better illustration (the original shape is (24, 32)). Each map in *x*-row and *y*-col of the figure represents the cosine similarities between the embedding vector in position (x, y) and the embedding vectors at other locations.

similarity values with the neighbour locations in 2D-grid and lower values with those far away from it. The results indicate the coarse 2D position information has been implicitly learned in the learnable position embedding. We suppose that the learning sources of the position information might be the 2D-structure groundtruth heatmaps and the similar features existing in the 1D-structure sequences. The model learns to build associations between position embedding and input sequences, as a result it can predict the target heatmaps with 2D Gaussian peaking at groundtruth keypoints locations.

In the paper, we find that position embedding helps to generalize better on unseen input resolutions, particularly 2D sine position embedding. We conjecture that 1) the models with a fixed receptive field may be hard to adapt the changes in scales; 2) building associations with *position information* encoded in Sine position embedding [9] may help model generalize better on different sizes.

## C. Transformer Encoder Layer

The Transformer Encoder layer [9] we used can be formulated as:

$$\begin{split} \mathbf{Z} &= \mathrm{LayerNorm}\left(\mathrm{MultiheadSelfAttention}\left(\mathbf{X}\right) + \mathbf{X}\right),\\ \mathbf{X}^{*} &= \mathrm{LayerNorm}\left(\mathrm{FFN}\left(\mathbf{Z}\right) + \mathbf{Z}\right), \end{split}$$

where X is the original input sequence that has not yet been added with position embedding. The position embedding will be added to X for computing querys and keys excluding values.  $X^*$  is the output sequence of the current Transformer Encoder layer, as the input sequence of next encoder layer. The formulations of Multihead Self-Attention and FFN are defined in [9].

#### **D.** Gradient Analysis

From the view of an activation at some location i of the predicted heatmaps, the network weights associating all input tokens across the whole image/sequence with this activation can be seen as a discriminator that judges the presence or absence of a certain keypoint at this location. As revealed by [5, 7, 1, 6], the gradient information can indicate the importance (sensitivity) of the input features to a specific output of a non-linear model. That assumption is based on that tiny change in the input (pixel/feature/token) with the most important feature value causes a large change in what the output of the model would be.

Suppose we have a trained model and a specific image,  $h_i \in \mathbb{R}^K$  is the scores for all K types of keypoints at location *i* of the predicted heatmaps;  $z_i \in \mathbb{R}^d$  is the intermediate feature outputted by the last self-attention laver before being fed into FFN. There is only a ReLU excluding the linear and convolutions<sup>1</sup> (head) layers after the last attention layer. ReLU (rectified linear unit) activation function in FFN can be empirically regarded as a negative contribution filter, which only retains positive contributions and maintains the linearity. Next, we choose numerator layout for computing the derivative of a vector with respect to a vector. We thus assume the mapping from  $z_i$  to  $h_i$ can be approximated as a linear function f with learned weights  $\mathbf{W}_{f} \in \mathbb{R}^{K \times d}$  and bias  $\mathbf{b} \in \mathbb{R}^{K}$  by computing the first-order Taylor expansion at a given local point  $z_i^0$ , *i.e.*,  $h_i \approx \mathbf{W}_f \mathbf{z}_i + \mathbf{b}, \ \mathbf{W}_f = \left. \frac{\partial h_i}{\partial \mathbf{z}_i} \right|_{\mathbf{z}_i^0}$ . Then we compute the partial derivative of  $h_i$  at location i of the output heatmaps w.r.t the token  $x_i$  at location j of the input sequence of the

last attention layer:

$$\begin{aligned} \frac{\partial \mathbf{h}_{i}}{\partial \mathbf{x}_{j}} &= \frac{\partial \mathbf{h}_{i}}{\partial \mathbf{z}_{i}} \frac{\partial \mathbf{z}_{i}}{\partial \mathbf{x}_{j}} \\ &= \frac{\partial f(\mathbf{z}_{i})}{\partial \mathbf{z}_{i}} (\mathbf{1} + \frac{\partial \mathbf{w}_{i} \mathbf{V}}{\partial \mathbf{x}_{j}}) \\ &\approx \mathbf{W}_{f} (\mathbf{1} + \frac{\partial \mathbf{w}_{i,0} \mathbf{v}_{0} + \ldots + \mathbf{w}_{i,j} \mathbf{v}_{j} + \ldots + \mathbf{w}_{i,L-1} \mathbf{v}_{L-1})}{\partial \mathbf{x}_{j}} \end{aligned}$$
$$\begin{aligned} &= \mathbf{W}_{f} (\mathbf{1} + \frac{\partial \mathbf{w}_{i,j} \mathbf{v}_{j}}{\partial \mathbf{x}_{j}}) \\ &= \mathbf{W}_{f} (\mathbf{1} + \frac{\partial \mathbf{A}_{i,j} \mathbf{W}_{v}^{\top} \mathbf{x}_{j}}{\partial \mathbf{x}_{j}}) \end{aligned}$$
(4)

where  $v_j \in \mathbb{R}^d$  is the value vector transformed by:  $v_j = \mathbf{W}_v^\top \boldsymbol{x}_j$ .  $\mathbf{A}_{i,j}$  is a scalar value that is computed by the dotproduct between  $q_i$  and  $k_j$ . We assume  $G := \frac{\partial h_i}{\partial x_j}$  as a function w.r.t. a given attention score  $\mathbf{A}_{i,j}$ . Under this assumption  $\mathbf{A}_{i,j}$  is deemed as an observed variable that has blocked its parent nodes. Then we define:

$$G \left( \mathbf{A}_{i,j} \right) = \mathbf{W}_{f} \left( \mathbf{1} + \frac{\partial \mathbf{A}_{i,j} \mathbf{W}_{v}^{\top} \boldsymbol{x}_{j}}{\partial \boldsymbol{x}_{j}} \right)$$

$$= \mathbf{W}_{f} \left( \mathbf{1} + \mathbf{A}_{i,j} \mathbf{W}_{v}^{\top} \right)$$

$$= \mathbf{A}_{i,j} \mathbf{W}_{f} \mathbf{W}_{v}^{\top} + \mathbf{W}_{f}$$

$$= \underbrace{\mathbf{A}_{i,j}}_{\text{Image-Specific: dynamic weights}} \cdot \underbrace{\mathbf{W}_{f} \cdot \mathbf{W}_{v}^{\top} + \mathbf{W}_{f}}_{\text{Learned: static weights}}$$

$$= \mathbf{A}_{i,j} \cdot \mathbf{K} + \mathbf{B}$$
(5)

where  $\mathbf{K}, \mathbf{B} \in \mathbb{R}^{K \times d}$  are static weights shared across all positions. We can see that the function *G* is approximately linear with  $\mathbf{A}_{i,j}$ , *i.e.*, the degrees of contribution to the prediction  $h_i$  directly depend on its attention scores at those locations.

The last attention layer in Transformer Encoder, whose attention scores are seen as the image-specific weights, aggregate contributions from all locations according to attention scores and finally form the maximum activations in the output heatmaps. Though the layers in FFN and head cannot be ignored<sup>2</sup>, they are position-wise operators, which almost linearly transform the attention scores from all the positions with the same transformation. In addition,  $\mathbf{Q} = (\mathbf{X} + \mathbf{P}) \mathbf{W}_q, \mathbf{K} = (\mathbf{X} + \mathbf{P}) \mathbf{W}_k, \mathbf{V} = \mathbf{X} \mathbf{W}_v$  where  $\mathbf{P}$  is the position embedding. Because  $\mathbf{A}_{i,j} \propto \mathbf{Q}_i \mathbf{K}_j^{\top}$ , the position embedding values also affect the attention scores to some extent.

 $<sup>^{1}</sup>a$  1  $\times$  1 convolution is also a position-wise linear layer; the 4  $\times$  4 deconvolution used in TP-R acts as the upsampling operation.

 $<sup>^{2}</sup>$ 1. Assuming that the used convolutions extract feature in a limited patch, the global interactions mostly occur at the attention layers. 2. The layer normalization does not affect the interactions between locations.



(a) **TP-R-A4:** predictions and dependency areas for input 1.



(c) TP-R-A4: predictions and dependency areas for input 2.



(e) TP-R-A4: predictions and dependency areas for input 3.



(g) TP-R-A4: predictions and dependency areas for input 4.



(i) TP-R-A4: predictions and dependency areas for input 5.



(k) TP-R-A4: predictions and dependency areas for input 5.

	AN CO			
Sec.			Ser .	Ye.

(b) **TP-H-A4:** predictions and dependency areas for **input 1**.



(d) TP-H-A4: predictions and dependency areas for input 2.



(f) **TP-H-A4:** predictions and dependency areas for input 3.



(h) TP-H-A4: predictions and dependency areas for input 4.



(j) TP-H-A4: predictions and dependency areas for input 5.



(1) **TP-H-A4:** predictions and dependency areas for input 5.

Figure 2. Predicted locations and the dependency areas for different types of keypoints in different models: TP-R-A4 (left column) and TP-H-A4 (right column). In each sub-figure, the first one is the original input image plotted with predicted skeleton. The other maps visualized by the defined dependency area  $(A_{i,:})$  of the attention matrix in the last layer with a threshold value (0.00075). The predicted location of a keypoint is annotated by a WHITE color pentagram (\*) in each sub-map. Redder area indicates higher attention scores.



(a) **TP-R-A4:** predictions and dependency areas of each keypoint in different attention layers.



(c) **TP-R-A4:** predictions and dependency areas of each keypoint in different attention layers.



(e) **TP-R-A4:** predictions and affect areas of each keypoint in different attention layers.

Atten. Layer 0	- 04	- Martin	E CONTRACTOR		10	100	(Pro-	48	1		3	1	1	-	Contraction of the second	R.		10.00	1. A.
Atten. Layer 1			200	200	200			<u> (</u>	1	3	000	10	200	383	300		1	120	1910
Atten. Layer 2	ġ,		5 120	<b>(</b> 10)	500			1	1	No.		1	0 lite	2	-	E.	Sec. 10		3.0%
Atten. Layer 3		(A)						100				130	-	0.0	38		Cire.	200	Ser.
r	lose	eve(l)	eve(r)	eariii	ear(r)	shp.(I)	sho.(r)	elb.(l)	elb.(r)	wri.(I)	wri.(r)	hipili	hip(r)	kne.(I)	kne.(r)	ank.(I)	ank.(r)	random	random

(b) **TP-H-A4:** predictions and dependency areas of each keypoint in different attention layers.

Atten. Layer 0	J.	đđ,	Sec.	200		B.	83	36	22	38	3	200	-85°	×.		38	<u> </u>	200	380	0.0
Atten. Layer 1	J.		and a	19.9°			ers?	300	30	est?	S.	Star Star	199	3	3.00					
Atten. Layer 2	J.	200	2.5	(1)) (1))	1. S	4.0	30		30		$(2^{12})$		100			<b>\$</b>		2.	100	10
Atten. Layer 3	j.	100°	20	C.S	35	85	er f	10	<b>3</b>	egges <sup>o</sup>	and a	340	ŝ	300	30	32	30	(). ().	des -	300
		nose	eye(l)	eve(r)	ear(I)	ear(r)	sho.(I)	sho.(r)	elb.(I)	elb.(r)	wri.(I)	wn.(r)	hipti)	hip(r)	kne.(I)	kne.(r)	ank.(I)	ank.(r)	randomi	andom

(d) **TP-H-A4:** predictions and dependency areas of each keypoint in different attention layers.



(f) **TP-H-A4:** predictions and affect areas of each keypoint in different attention layers.

Figure 3. Dependency areas (the first two rows) and Affected areas (the last row) in different attention layers for different input images.

Backbone	ResNet-S
Stem	Conv-k7-s2-c64, BN, ReLU Pooling-k3-s2
Blocks	3×Bottleneck-c64 Bottleneck-s2-c128 3×Bottleneck-c128 Conv-k1-s1-c256

Table 1. The detailed configurations for ResNet-S. Conv-k7-s2c64 means a convolutional layer with  $7 \times 7$  kernel size, 2 stride, and 64 output channels, followed by a BN and ReLU; the same below. The Bottleneck-c64 includes Conv-k1-s1-c64-BN-ReLU, Conv-k3-s1-c64-BN-ReLU, and Conv-k1-s1-c256-BN. Bottleneck-c128 includes Conv-k1-s1-c128-BN-ReLU, Conv-k3s1-c128-BN-ReLU, and Conv-k1-s1-c512-BN. See details in [3].

Backbone	HRNet-S-W32(48)
Stem	Conv-k3-s2-c64, BN, ReLU Conv-k3-s2-c64, BN, ReLU 4×Bottleneck-c64
Blocks	transition1~stage2 transition2~stage3 Conv-k1-s1-c64(92)

Table 2. The detailed configurations for HRNet-S-W32(48). More detailed information about the transition layer and stage blocks are described in the HRNet paper [8].

# **E.** Architecture Details

We report the architecture details of ResNet-S and HRNet-S-W32(48) in Tab. 1 and Tab. 2. The ResNet-S\* only differs from ResNet-S in that ResNet-S\* has 10 Bottleneck-c128 blocks. More details about HRNet-W32 and HRNet-W48 are described in [8].

## F. More Attention Maps Visualizations

In this section, we show more visualization results of the attention maps from TP-R-A4 (TransPose-R-A4) and TP-H-A4 (TransPose-H-A4) models. The attention maps of the last attention layers of two models are shown in Fig. 2. The attention maps in different attention layers of two models are shown in Fig. 3.

## References

- Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller, and Wojciech Samek. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PloS one*, 10(7):e0130140, 2015. 2
- [2] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-toend object detection with transformers. In *ECCV*, pages 213– 229, Cham, 2020. 1

- [3] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016. 4
- [4] Niki Parmar, Ashish Vaswani, Jakob Uszkoreit, Lukasz Kaiser, Noam Shazeer, Alexander Ku, and Dustin Tran. Image transformer. In *ICML*, 2018. 1
- [5] Wojciech Samek, Grégoire Montavon, Andrea Vedaldi, Lars Kai Hansen, and Klaus-Robert Müller. *Explainable AI: interpreting, explaining and visualizing deep learning*, volume 11700. Springer Nature, 2019. 2
- [6] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Gradcam: Visual explanations from deep networks via gradientbased localization. In *ICCV*, pages 618–626, 2017. 2
- [7] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. In *ICLR (Workshop Poster)*, 2013. 2
- [8] Ke Sun, Bin Xiao, Dong Liu, and Jingdong Wang. Deep highresolution representation learning for human pose estimation. In *CVPR*, June 2019. 4
- [9] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, pages 5998–6008, 2017. 1, 2