

# End-to-End Transformer Based Model for Image Captioning (Supplementary Materials)

Anonymous submission: 8053

Anonymous submission

## Supplementary Materials

### More Experimental Details

**Source Code** The source code for the construction of our PureT is provided as supplementary material, which can be found in folder `PureT_source_code(Excerpt)/models`. We will release trained models and the full source code.

**Experimental Results on MSCOCO** The detailed captions generated by our models reported in Table 1 and Table 2 of the main paper are provided as supplementary material, organized in the following format:

```
[  
 {  
   "image_id": 391895,  
   "caption": "a man walking a dirt bike on a dirt road"  
 },  
 ...  
 ]
```

which include:

- Offline test results of single model:  
`MSCOCO_results/Singlemodel_offline_results`
- Offline test results of ensemble model:  
`MSCOCO_results/Ensemblemodel_offline_results`
- Online test results of ensemble model:  
`MSCOCO_results/Ensemblemodel_online_results`

**Experimental Environments** We use Intel Xeon E5-2698 v4 2.2 GHz (20-Core) with Tesla V100 for our experiments. The operating system is Ubuntu Desktop Linux OS 18.04, and the main software environments include Python 3.7.4 and PyTorch 1.5.1.

**Visualization of Attention Weights** In order to facilitate the intuitive observation of the transformation of attention area during the caption generation process, we visualize the attention weights of Cross MSA Module in decoder. Specifically, we visualize the Cross MSA module in the last block of decoder at each timestep. The size of attention weights is [8, 12, 12], where 8 is the number of heads, we calculate the average values of the 8 heads to obtain the final attention weights  $\alpha$  with the size of [12, 12] for visualization. We only

retain the top-20 weight values and set the remaining values to 0 for better visualization.

Furthermore, we need to scale  $\alpha$  from [12, 12] to the original size [H, W] of the input image to obtain a brightness mask. This can be implemented by applying the `resize(·)` and `pyramid_expand(·)` methods of `skimage.transform` in Python.

### Additional Ablation Study

To quantify the influence of different features extracted by different backbone models, we adopt different image captioning models, as baseline models and ablate them with different configurations of backbone models as shown in Table 6. The baseline models include:  $\mathcal{M}^2$  Transformer (Cornia et al. 2020), X-Transformer (Pan et al. 2020) and standard Transformer (Vaswani et al. 2017). The backbone models include: Faster R-CNN (Ren et al. 2017) in conjunction with ResNet-101, which is adopted in (Anderson et al. 2018); Faster R-CNN in conjunction with ResNeXt-101, which is adopted in (Jiang et al. 2020); SwinTransformer (Liu et al. 2021).

As we can see, grid features extracted by SwinTransformer can achieve significant performance improvement compared with region features extracted by ResNet-101 and grid features extracted by ResNeXt-101.

In terms of  $\mathcal{M}^2$  Transformer and X-Transformer, the backbone models of ResNet-101 and ResNeXt-101 have similar performance. The backbone model of SwinTransformer comprehensively improves scores of all metrics, which boosts the CIDEr score more than 3.7% in  $\mathcal{M}^2$  Transformer especially. Note that the backbone model with  $N = 3$  has a better performance than  $N = 6$  in X-Transformer, which indicates the superiority of SwinTransformer in image captioning and allows us to explore more tiny and efficient models and apply it in more actual scenes. In terms of standard Transformer, the backbone model of SwinTransformer reaches an excellent performance and is even better than  $\mathcal{M}^2$  Transformer and X-Transformer in scores of METEOR, CIDEr and SPICE. In terms of our PureT, the backbone of SwinTransformer also achieves a better performance than ResNeXt-101.

In general, in our extensive experiments, we find that the backbone models of CNN (e.g. Faster RCNN in conjunction with ResNet-101 or ResNeXt-101) are more suitable for us-

Baseline Models	Backbone	Feat. Type	Feat. Size	N	B-1	B-2	B-3	B-4	M	R	C	S
$\mathcal{M}^2$ Transformer	ResNet-101	Region	(10–100)	3 <sup>†</sup>	80.8	-	-	39.1	29.2	58.6	131.2	22.6
	ResNeXt-101	Grid	7 × 7	3 <sup>‡</sup>	80.8	-	-	38.9	29.1	58.5	131.7	22.6
	SwinTransformer	Grid	12 × 12	3	81.8	66.8	52.6	40.5	29.6	59.9	135.4	23.3
X-Transformer	ResNet-101	Region	(10–100)	6 <sup>†</sup>	80.9	65.8	51.5	39.7	29.5	59.1	132.8	23.4
	ResNeXt-101	Grid	7 × 7	6 <sup>‡</sup>	81.0	-	-	39.7	29.4	58.9	132.5	23.1
	SwinTransformer	Grid	12 × 12	6	81.4	66.3	52.0	39.9	29.5	59.5	133.7	23.4
	SwinTransformer	Grid	12 × 12	3	81.9	66.7	52.3	40.1	29.6	59.6	134.8	23.4
standard Transformer	ResNet-101	Region	(10–100)	3	80.0	64.9	50.5	38.7	29.0	58.6	130.1	22.9
	ResNeXt-101	Grid	7 × 7	3 <sup>‡</sup>	81.2	-	-	39.0	29.2	58.9	131.7	22.6
	ResNeXt-101	Grid	12 × 12	3	80.8	65.8	51.4	39.4	29.4	59.2	132.8	23.2
	SwinTransformer	Grid	12 × 12	3	81.6	66.5	52.0	39.8	29.9	59.6	136.4	23.8
PureT	ResNeXt-101	Grid	12 × 12	3	80.7	65.9	51.7	39.9	29.2	59.1	131.8	23.0
	SwinTransformer	Grid	12 × 12	3	<b>82.1</b>	<b>67.3</b>	52.0	<b>40.9</b>	<b>30.2</b>	<b>60.1</b>	<b>138.2</b>	<b>24.2</b>

Table 6: Performance comparison of different configuration of backbone models. ResNet-101 and ResNeXt-101 indicate Faster R-CNN in conjunction with them respectively. Region features extracted by ResNet-101 have adaptive size of 10 to 100. Grid features extracted by ResNeXt-101 can be extracted in the size of 12 × 12 or 7 × 7 by average pooling as need. Grid features (SwinTransformer) are extracted in the size of 12 × 12. N denotes the number of encoder and decoder blocks, superscript † indicates that the results are from the respectively official paper and ‡ indicates that the results are from (Luo et al. 2021), and other results come from our experiments.

ing LSTM or Transformer with non-standard MSA (e.g. X-Transformer) as decoder, and the backbone of SwinTransformer is more suitable for using Transformer with standard MSA (e.g.  $\mathcal{M}^2$  Transformer, standard Transformer and our PureT) as decoder. Therefore, we intend to explore a lighter and simpler Transformer-based model in our future work.

### Additional Visualization Examples

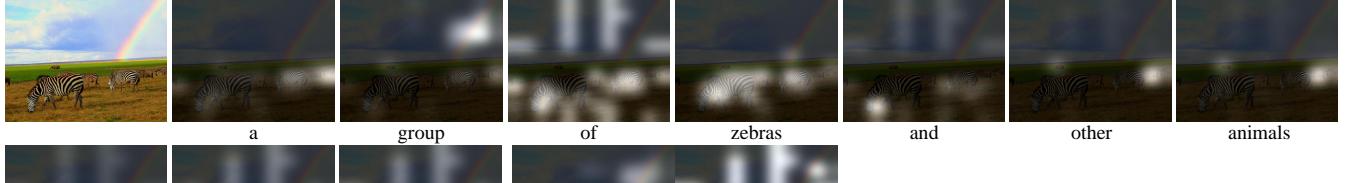
Figure 5 gives a comparison of three extra visualization examples of the standard Transformer and our PureT. It can be seen that both standard Transformer and our PureT can pay attention to the correct area when generating words. But our PureT can catch more fine-grained details and generate more accurate and descriptive captions.

Furthermore, we illustrate more examples of captions generated by standard Transformer,  $\mathcal{M}^2$  Transformer and our PureT, as shown in Figure 6, Figure 7 and Figure 8. The red marked words indicate wrong information or incomplete sentence, and the green marked words indicate more fine-grained details. For example,  $\mathcal{M}^2$  Transformer generates an incomplete sentence “a boat in the water with a mountain in the” for the 4-th image from left of Figure 6, while our PureT generates “a small boat in a large lake with mountains in the background” and contains more detailed adjectives “small” and “large”;  $\mathcal{M}^2$  Transformer generates an incorrect sentence “a group of men playing tennis on a field” for the 6-th image from right of Figure 6, in which the splicing of “three pictures” is incorrectly recognized as “a group of people”. Our model generates a correct sentence “three pictures of men playing tennis in a field”. For another example in the 8-th image from right of Figure 7,  $\mathcal{M}^2$  Transformer generates “a young boy standing in front of a refrigerator” and incorrectly identifies the “door” as “refrigerator”, standard Transformer generates a simple and correct sentence “a young boy standing in front of a door”, and our PureT generates “a young boy wearing a tie standing in front of a door”

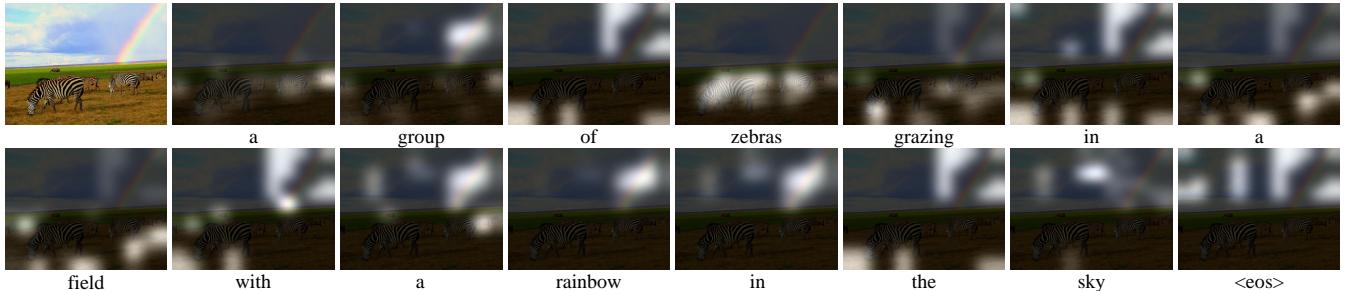
with more details information “wearing a tie”.

### References

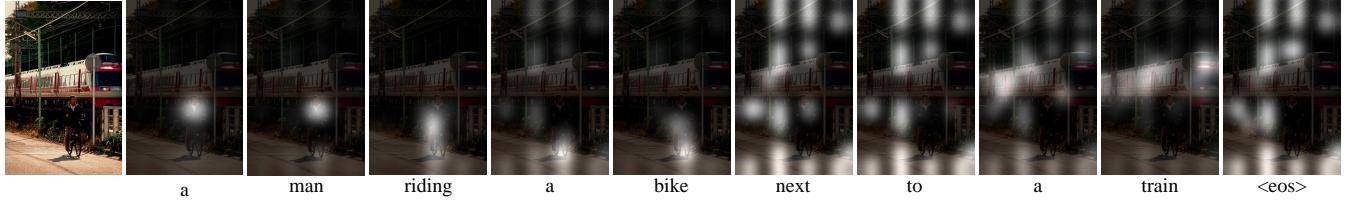
- Anderson, P.; He, X.; Buehler, C.; Teney, D.; Johnson, M.; Gould, S.; and Zhang, L. 2018. Bottom-Up and Top-Down Attention for Image Captioning and Visual Question Answering. In *Proceedings of the CVPR*, 6077–6086.
- Cornia, M.; Stefanini, M.; Baraldi, L.; and Cucchiara, R. 2020. Meshed-Memory Transformer for Image Captioning. In *Proceedings of the CVPR*, 10575–10584.
- Jiang, H.; Misra, I.; Rohrbach, M.; Learned-Miller, E. G.; and Chen, X. 2020. In Defense of Grid Features for Visual Question Answering. In *Proceedings of the CVPR*, 10264–10273. IEEE.
- Liu, Z.; Lin, Y.; Cao, Y.; Hu, H.; Wei, Y.; Zhang, Z.; Lin, S.; and Guo, B. 2021. Swin Transformer: Hierarchical Vision Transformer using Shifted Windows. *arXiv preprint arXiv:2103.14030*.
- Luo, Y.; Ji, J.; Sun, X.; Cao, L.; Wu, Y.; Huang, F.; Lin, C.; and Ji, R. 2021. Dual-level Collaborative Transformer for Image Captioning. In *Proceedings of the AAAI*, 2286–2293. AAAI Press.
- Pan, Y.; Yao, T.; Li, Y.; and Mei, T. 2020. X-Linear Attention Networks for Image Captioning. In *Proceedings of the CVPR*, 10968–10977.
- Ren, S.; He, K.; Girshick, R. B.; and Sun, J. 2017. Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks. *IEEE Trans. Pattern Anal. Mach. Intell.*, 39(6): 1137–1149.
- Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A. N.; Kaiser, L.; and Polosukhin, I. 2017. Attention is All you Need. In *Proceedings of the NIPS*, 5998–6008.



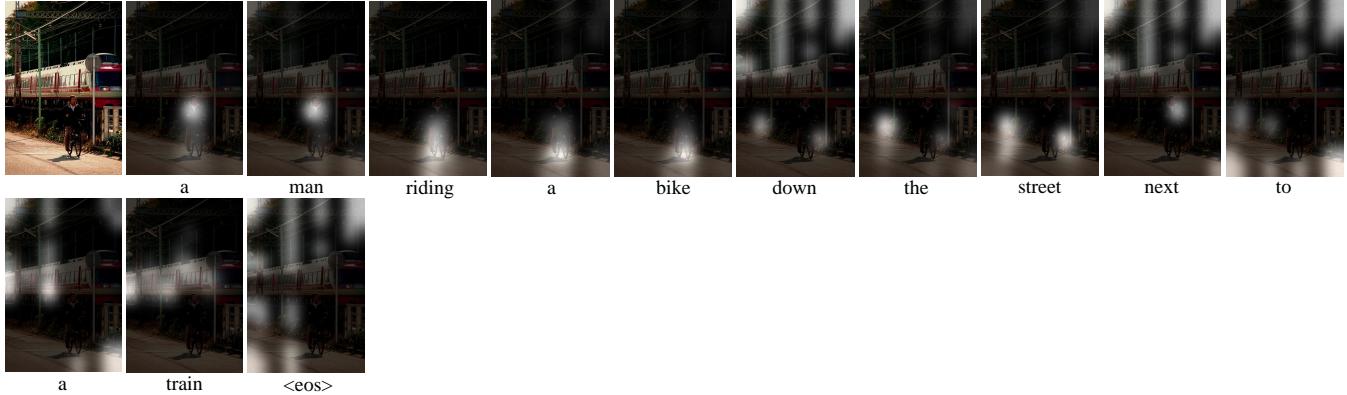
Standard Transformer: a group of zebras and other animals grazing in a field



Our: a group of zebras grazing in a field with a rainbow in the sky.



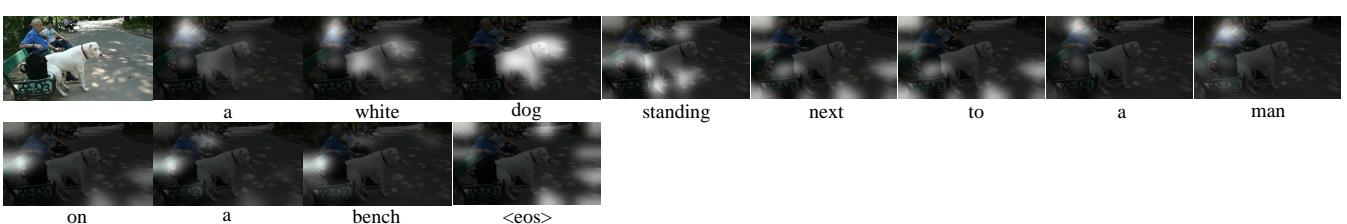
Standard Transformer: a man riding a bike next to a train.



Our: a man riding a bike down the street next to a train.



Standard Transformer: a white dog sitting on a bench.



Our: a white dog standing next to a man on a bench.

Figure 5: Visualization of attention heatmap during caption generation process of standard Transformer and our PureT.



Figure 6: More examples of captions generated by standard Transformer,  $\mathcal{M}^2$  Transformer and our PureT. The red marked words indicate wrong information or incomplete sentence, and the green marked words indicate more fine-grained details.

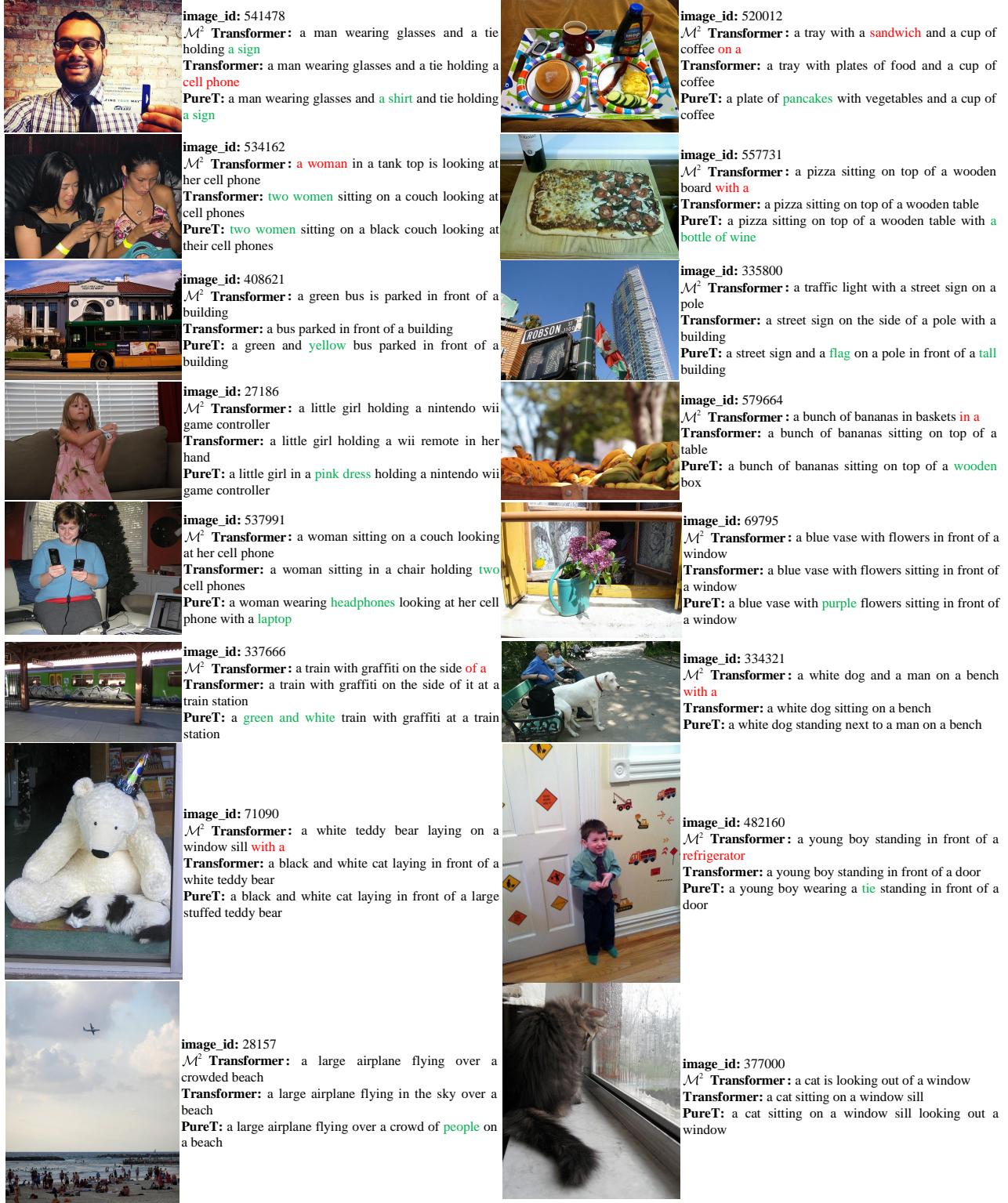


Figure 7: More examples of captions generated by standard Transformer,  $\mathcal{M}^2$  Transformer and our PureT. The red marked words indicate wrong information or incomplete sentence, and the green marked words indicate more fine-grained details.



image\_id: 360318

$\mathcal{M}^2$  Transformer: a brown bear sitting on a rock eating food  
Transformer: a monkey eating a piece of food on a rock  
PureT: a monkey sitting on a rock eating food



image\_id: 442993

$\mathcal{M}^2$  Transformer: a building with a clock tower on top of a  
Transformer: a building with a clock tower on top of it  
PureT: a building with a clock tower in front of the water



image\_id: 469618

$\mathcal{M}^2$  Transformer: a group of people sitting around a table with plates of food  
Transformer: a group of people sitting at a table with plates of food  
PureT: a group of people sitting at a picnic table with plates of food



image\_id: 236068

$\mathcal{M}^2$  Transformer: a man in a top hat talking on a cell phone  
Transformer: a man in a suit talking on a cell phone  
PureT: a man wearing a top hat talking on a cell phone next to an old car



image\_id: 444302

$\mathcal{M}^2$  Transformer: a street sign on the side of a road  
Transformer: a street sign on the side of a road  
PureT: a black and white photo of a street sign on the side of a road



image\_id: 483108

$\mathcal{M}^2$  Transformer: a man riding a bike with a train in the  
Transformer: a man riding a bike next to a train  
PureT: a man riding a bike down the street next to a train

Figure 8: More examples of captions generated by standard Transformer,  $\mathcal{M}^2$  Transformer and our PureT. The red marked words indicate wrong information or incomplete sentence, and the green marked words indicate more fine-grained details.