

# Supplementary Material - CNN Attention Guidance for Improved Orthopedics Radiographic Fracture Classification

Zhibin Liao, Kewen Liao, Haifeng Shen, Marouska F. van Boxel, Jasper Prijs, Ruurd L. Jaarsma, Job N. Doornberg, Anton van den Hengel, Johan W. Verjans

In this supplementary material, we show the experiment results of the proposed method on the ResNet-101 [1] and various attentive backbones that demonstrate a consistent performance across the experimented models.

## I. RESNET-101

The results using the ResNet-101 backbone are displayed in Table I. In summary, we found the baseline performances (*i.e.*, classification models without attention guidance) are improved from the respective ResNet-50 [1] counterparts (scaphoid (s-512) study-wise is up from  $76.4\% \pm 0.9\%$  to  $79.6\% \pm 1.5\%$ ,  $p < 0.01$ , ankle (Test Set) from  $69.0\% \pm 2.8\%$  to  $73.4\% \pm 4.0\%$ ,  $p = 0.08$ ), showing that the larger/deeper neural network is able to generalize better. On the other hand, the attention guidance models are still able to provide consistent improvements over the ResNet-101 baseline results, which demonstrates the effectiveness of the proposed method. In addition, the best results collected from the ResNet-101 backbone models (scaphoid  $85.8\% \pm 1.3\%$  using Scribble at  $\lambda = 0.1$ , and ankle  $79.6\% \pm 1.1\%$  using Segm. at  $\lambda = 1$ ) are not significantly better than those from the paper's ResNet-50 models (scaphoid  $84.2\% \pm 0.8\%$  using BBOX at  $\lambda = 1$ ,  $p = 0.06$ , and ankle  $80.6\% \pm 3.8\%$  using Segm. at  $\lambda = 0.1$ ,  $p = 0.60$ ).

TABLE I

THE SCAPHOID AND ANKLE FRACTURE EXPERIMENT RESULTS IN CLASSIFICATION ACCURACY (%), USING THE RESNET-101 BACKBONE WITH  $512 \times 512$  INPUT SIZE.

Methods	Scaphoid					
	Image-wise			Study-wise		
Baseline	$76.9 \pm 1.0$			$79.6 \pm 1.5$		
Attn. Guid.	Scribble	BBOX	Segm.	Scribble	BBOX	Segm.
$\lambda = 0.01$	$78.2 \pm 1.2$	$77.0 \pm 2.2$	$77.3 \pm 1.5$	$83.6 \pm 1.3$	$81.6 \pm 3.4$	$81.6 \pm 2.6$
$\lambda = 0.1$	$79.7 \pm 1.1$	<b><math>78.5 \pm 0.8</math></b>	<b><math>79.0 \pm 1.6</math></b>	<b><math>85.8 \pm 1.3</math></b>	$82.8 \pm 1.6$	$83.4 \pm 2.6$
$\lambda = 1$	<b><math>79.9 \pm 1.3</math></b>	$78.3 \pm 1.0$	$78.7 \pm 1.0$	$84.8 \pm 1.9$	<b><math>83.0 \pm 1.4</math></b>	<b><math>85.0 \pm 1.2</math></b>
$\lambda = 5$	$73.7 \pm 1.0$	$77.9 \pm 1.1$	$77.5 \pm 1.1$	$78.4 \pm 2.2$	$82.2 \pm 0.8$	$82.4 \pm 2.0$
Methods	Ankle					
	Image-wise			Study-wise		
Baseline	$69.1 \pm 3.1$			$73.4 \pm 4.0$		
Attn. Guid.	Scribble	BBOX	Segm.	Scribble	BBOX	Segm.
$\lambda = 0.1$	$71.0 \pm 1.9$	$71.2 \pm 0.8$	$69.2 \pm 1.5$	$76.8 \pm 2.3$	$76.6 \pm 2.5$	$75.0 \pm 1.7$
$\lambda = 0.5$	<b><math>71.9 \pm 1.8</math></b>	<b><math>72.3 \pm 1.6</math></b>	$70.5 \pm 2.9$	<b><math>77.4 \pm 1.1</math></b>	<b><math>79.2 \pm 1.5</math></b>	$76.0 \pm 4.6$
$\lambda = 1$	$70.2 \pm 2.7$	$71.3 \pm 2.5$	<b><math>73.6 \pm 0.7</math></b>	$76.2 \pm 1.8$	$77.6 \pm 1.3$	<b><math>79.6 \pm 1.1</math></b>

## II. ATTENTIVE BACKBONES

The results of all three types of human guidance (*i.e.*, scribble, bounding box, and segmentation) on the attentive

backbones (*i.e.*, residual attention [2], self attention [3], and pyramid attention [4]) can be found in Table II. We use the ResNet-50 backbone for these attentive models and set input size to 512. The  $\lambda$  values are set as 1 for the scaphoid experiments and 0.5 for the fracture experiments as they were the optimal values from the previous experiments (*i.e.*, models backbone by the non-attentive vanilla ResNet-50) in the main paper.

In summary, these experiment results convey the similar observation, *i.e.*, with the use of human guidance, the classification results are much improved. On the scaphoid (s-512) dataset, the best classification performance is observed with the combination of self attention + scribble type guidance ( $84.6\% \pm 1.7\%$ ), on par with the previous optimal result ( $84.2\% \pm 0.8\%$ , vanilla ResNet50 + BBOX). On the ankle test set, the optimal result is  $81.2\% \pm 1.3\%$  (self attention + Segm.), which is also on par with the previous optimal result ( $80.4\% \pm 3.9\%$ , vanilla ResNet50 + BBOX). The only exception is with the pyramid attention on the ankle test set, where the model performs significantly worse both with and without human guidance. This may be a result of the backbone model in which the addition of pyramid attention mechanism significantly changed the energy landscape of the optimization on the ankle dataset.

## REFERENCES

- [1] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [2] F. Wang, M. Jiang, C. Qian, S. Yang, C. Li, H. Zhang, *et al.*, "Residual attention network for image classification," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 3156–3164.
- [3] H. Zhang, I. Goodfellow, D. Metaxas, and A. Odena, "Self-attention generative adversarial networks," in *International conference on machine learning*, PMLR, 2019, pp. 7354–7363.
- [4] H. Li, P. Xiong, J. An, and L. Wang, "Pyramid attention network for semantic segmentation," *arXiv preprint arXiv:1805.10180*, 2018.

TABLE II

THE HUMAN-GUIDED ATTENTION EXPERIMENT RESULTS WITH THE USE OF THE ATTENTIVE BACKBONES, *i.e.*, RESIDUAL ATTENTION [47], SELF ATTENTION [48], AND PYRAMID ATTENTION [49]. THE NO GUIDANCE ENTRY REPEATS THE RESULTS FOR THE RESPECTIVE MODELS IN TABLE III AND IV IN THE MAIN MANUSCRIPT FOR CONVENIENCE.

Dataset Backbone	Scaphoid						Ankle					
	Residual Attn.		Self Attn.		Pyramid Attn.		Residual Attn.		Self Attn.		Pyramid Attn.	
	Image	Study	Image	Study	Image	Study	Image	Study	Image	Study	Image	Study
No Guid.	75.3 $\pm$ 3.4	78.8 $\pm$ 2.4	75.5 $\pm$ 1.7	79.2 $\pm$ 1.8	75.7 $\pm$ 0.6	79.2 $\pm$ 2.7	62.5 $\pm$ 1.6	65.0 $\pm$ 1.2	68.9 $\pm$ 2.5	74.6 $\pm$ 3.1	<b>55.4 <math>\pm</math> 6.2</b>	<b>59.4 <math>\pm</math> 9.6</b>
Scribble	77.3 $\pm$ 0.8	82.6 $\pm$ 2.7	<b>78.7 <math>\pm</math> 1.0</b>	<b>84.6 <math>\pm</math> 1.7</b>	77.3 $\pm$ 1.0	81.0 $\pm$ 1.4	72.7 $\pm$ 2.4	<b>79.8 <math>\pm</math> 1.6</b>	71.9 $\pm$ 3.0	78.0 $\pm$ 4.4	50.7 $\pm$ 2.3	51.4 $\pm$ 1.5
BBOX	<b>78.9 <math>\pm</math> 1.0</b>	82.0 $\pm$ 2.4	77.6 $\pm$ 1.4	83.6 $\pm$ 1.5	77.7 $\pm$ 0.9	82.2 $\pm$ 1.8	71.3 $\pm$ 1.3	78.0 $\pm$ 1.6	72.3 $\pm$ 2.1	77.8 $\pm$ 2.2	53.2 $\pm$ 4.0	55.0 $\pm$ 5.4
Segm.	77.6 $\pm$ 1.8	<b>83.4 <math>\pm</math> 0.6</b>	77.2 $\pm$ 1.9	83.0 $\pm$ 3.1	<b>79.4 <math>\pm</math> 1.0</b>	<b>83.4 <math>\pm</math> 1.8</b>	<b>73.2 <math>\pm</math> 1.9</b>	79.2 $\pm$ 1.9	<b>74.0 <math>\pm</math> 2.1</b>	<b>81.2 <math>\pm</math> 1.3</b>	53.4 $\pm$ 5.0	55.2 $\pm$ 4.8