

# BiViT: Extremely Compressed Binary Vision Transformers

## Supplementary Material

Yefei He<sup>1</sup>   Zhenyu Lou<sup>1</sup>   Luoming Zhang<sup>1</sup>   Jing Liu<sup>2</sup>   Weijia Wu<sup>1</sup>  
Hong Zhou<sup>1†</sup>   Bohan Zhuang<sup>2†</sup>

<sup>1</sup>Zhejiang University, China

<sup>2</sup>ZIP Lab, Monash University, Australia

### A. More experimental results on COCO

In this section, we present additional experimental results of BiViT on the COCO 2017 validation set. As shown in Table A, our method outperforms the competitive BiBERT by a large margin in all metrics.

Table A. Comparisons of different methods and backbones on COCO 2017 validation set.

Framework	Method	Task	AP	AP <sub>50</sub>	AP <sub>75</sub>	AP <sub>S</sub>	AP <sub>M</sub>	AP <sub>L</sub>
Mask R-CNN	FP	Object Detection	43.7	66.6	47.7	28.5	47.0	57.3
	BiBERT		32.0	53.9	33.7	19.6	34.1	41.3
	Ours		<b>33.4</b>	<b>55.0</b>	<b>35.2</b>	<b>20.8</b>	<b>35.6</b>	<b>42.0</b>
	FP	Instance Segmentation	39.8	63.3	42.7	24.2	43.1	54.6
	BiBERT		30.4	51.0	31.9	16.7	32.5	41.4
	Ours		<b>31.6</b>	<b>51.7</b>	<b>33.4</b>	<b>18.3</b>	<b>34.0</b>	<b>42.1</b>
Cascade Mask R-CNN	FP	Object Detection	48.1	67.1	52.2	30.4	51.5	63.1
	BiBERT		39.2	57.3	42.5	23.7	41.3	51.1
	Ours		<b>40.8</b>	<b>59.2</b>	<b>44.1</b>	<b>25.2</b>	<b>43.3</b>	<b>52.8</b>
	FP	Instance Segmentation	41.7	64.4	45.0	41.7	64.4	45.0
	BiBERT		34.5	54.5	36.8	19.1	36.5	46.6
	Ours		<b>35.7</b>	<b>56.5</b>	<b>38.2</b>	<b>20.1</b>	<b>38.2</b>	<b>48.5</b>

<sup>†</sup>H. Zhou and B. Zhuang are corresponding authors.