

Model	#Trainable Params	Flickr30K Zero-shot (1K test set)						COCO Fine-tuned (5K test set)					
		Image $\rightarrow$ Text			Text $\rightarrow$ Image			Image $\rightarrow$ Text			Text $\rightarrow$ Image		
		R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10
<i>Dual-encoder models</i>													
CLIP (Radford et al., 2021)	428M	88.0	98.7	99.4	68.7	90.6	95.2	-	-	-	-	-	-
ALIGN (Jia et al., 2021)	820M	88.6	98.7	99.7	75.7	93.8	96.8	77.0	93.5	96.9	59.9	83.3	89.8
FILIP (Yao et al., 2022)	417M	89.8	99.2	99.8	75.0	93.4	96.3	78.9	94.4	97.4	61.2	84.3	90.6
Florence (Yuan et al., 2021)	893M	90.9	99.1	-	76.7	93.6	-	81.8	95.2	-	63.2	85.7	-
BEIT-3(Wang et al., 2022b)	1.9B	94.9	99.9	<b>100.0</b>	81.5	95.6	97.8	<u>84.8</u>	<u>96.5</u>	<u>98.3</u>	<u>67.2</u>	<b>87.7</b>	<b>92.8</b>
<i>Fusion-encoder models</i>													
UNITER (Chen et al., 2020)	303M	83.6	95.7	97.7	68.7	89.2	93.9	65.7	88.6	93.8	52.9	79.9	88.0
OSCAR (Li et al., 2020)	345M	-	-	-	-	-	-	70.0	91.1	95.5	54.0	80.8	88.5
VinVL (Zhang et al., 2021)	345M	-	-	-	-	-	-	75.4	92.9	96.2	58.8	83.5	90.3
<i>Dual encoder + Fusion encoder reranking</i>													
ALBEF (Li et al., 2021)	233M	94.1	99.5	99.7	82.8	96.3	98.1	77.6	94.3	97.2	60.7	84.3	90.5
BLIP (Li et al., 2022)	446M	96.7	<b>100.0</b>	<b>100.0</b>	86.7	97.3	98.7	82.4	95.4	97.9	65.1	86.3	91.8
<b>BLIP-2</b> ViT-L	474M	<u>96.9</u>	<b>100.0</b>	<b>100.0</b>	<u>88.6</u>	<u>97.6</u>	<b>98.9</b>	83.5	96.0	98.0	66.3	86.5	91.8
<b>BLIP-2</b> ViT-g	1.2B	<b>97.6</b>	<b>100.0</b>	<b>100.0</b>	<b>89.7</b>	<b>98.1</b>	<b>98.9</b>	<b>85.4</b>	<b>97.0</b>	<b>98.5</b>	<b>68.3</b>	<b>87.7</b>	<u>92.6</u>