

Model	# Pre-train	Use VG-QA	Test-Dev		Test-Standard		
	Images	For Fine-tuning	Overall	Yes/No	Number	Other	Overall
ViLT _{BASE} [29]	4M	✗	71.26	-	-	-	-
UNITER _{BASE} [5]	4M	✓	72.70	-	-	-	72.91
VILLA _{BASE} [16]	4M	✓	73.59	-	-	-	73.67
UNIMO _{BASE} [37]	4M	✗	73.79	-	-	-	74.02
ALBEF _{BASE} [35]	4M	✓	74.54	-	-	-	74.70
ALBEF _{BASE} [35]	14M	✓	75.84	-	-	-	76.04
VinVL _{BASE} [91]	5.7M	✗	75.95	-	-	-	76.12
VLMO _{BASE} [77]	4M	✗	76.64	-	-	-	76.89
BLIP _{BASE} [34]	14M	✓	77.54	-	-	-	77.62
METER-CLIP-ViT _{BASE} [13]	4M	✗	77.68	92.49	58.07	69.20	77.64
OFA _{BASE} [75]	54M	✓	77.98	-	-	-	78.07
SimVLM _{BASE} [78]	1.8B	✗	77.87	-	-	-	78.14
BLIP _{BASE} [34]	129M	✓	78.24	-	-	-	78.17
BLIP _{BASE} -CapFilt-L [34]	129M	✓	78.25	-	-	-	78.32
BRIDGE-TOWER _{BASE} (Ours)	4M	✗	78.66	92.92	60.69	70.51	78.73
BRIDGE-TOWER _{BASE} (Ours)	4M	✓	79.10	93.06	62.19	70.69	79.04
UNITER _{LARGE} [5]	4M	✓	73.82	-	-	-	74.02
VILLA _{LARGE} [16]	4M	✓	74.69	-	-	-	74.87
UNIMO _{LARGE} [37]	4M	✗	75.06	-	-	-	75.27
VinVL _{LARGE} [91]	5.7M	✗	76.52	92.04	61.50	66.68	76.63
SimVLM _{LARGE} [78]	1.8B	✗	79.32	-	-	-	79.56
VLMO _{LARGE} [77]	4M	✗	79.94	-	-	-	79.98
SimVLM _{HUGE} [78]	1.8B	✗	80.03	93.29	66.54	72.23	80.34
METER-CoSwin _{HUGE} [13]	14M	✗	80.33	94.25	64.37	72.30	80.54
OFA _{LARGE} [75]	54M	✓	80.43	93.32	67.31	72.71	80.67
BRIDGE-TOWER _{LARGE} (Ours)	4M	✗	81.25	94.69	64.58	73.16	81.15
BRIDGE-TOWER _{LARGE} (Ours)	4M	✓	81.52	94.80	66.01	73.45	81.49