Image-to-Text Recall@1 per Layer, Mask with Nearest Normal Token 55 50 (%) Recall@1 (Mask MA Tokens 40 Mask Artifact Tokens Mask MA & Artifact Tokens 35 Mask Normal Tokens = #MA Mask Normal Tokens = #Artifact Mask Normal Tokens = #MA & Artifact 30 Baseline (I2T): 51.44% 10 20 Layer