

	R Precision (top 1)↑	R Precision (top 2)↑	R Precision (top 3)↑	FID ↓	Diversity →	Multimodality ↑
Real	$0.4541 \pm 0.0060$	$0.6612 \pm 0.0065$	$0.7704 \pm 0.0054$	$0.0010 \pm 0.0001$	$9.1322 \pm 0.0943$	–
MDM-Original	$0.3883 \pm 0.0158$	$0.5828 \pm 0.0091$	$0.6945 \pm 0.0026$	$0.5243 \pm 0.1885$	$9.3325 \pm 0.0594$	$2.6270 \pm 0.0539$
MDM-1	$0.4495 \pm 0.0035$	$0.6576 \pm 0.0060$	$0.7681 \pm 0.0053$	$0.6466 \pm 0.0774$	$9.7104 \pm 0.1696$	$2.5973 \pm 0.1703$
MDM-2	$0.4546 \pm 0.0029$	$0.6606 \pm 0.0033$	$0.7681 \pm 0.0040$	$0.6042 \pm 0.0868$	$9.0701 \pm 0.1454$	$2.4932 \pm 0.0947$
MDM-3	$0.3979 \pm 0.0127$	$0.5877 \pm 0.0123$	$0.6992 \pm 0.0052$	$0.5602 \pm 0.0817$	$9.3052 \pm 0.1416$	$2.5170 \pm 0.0678$

1. DiT Block with adaLN-Zero
2. DiT Block with Cross-Attention
3. DiT Block with In-Context Conditioning

	R Precision (top 1)↑	R Precision (top 2)↑	R Precision (top 3)↑	FID ↓	Diversity →	Multimodality ↑
SWA-1	$0.4102 \pm 0.0081$	$0.6059 \pm 0.0083$	$0.7205 \pm 0.0087$	$1.1091 \pm 0.1846$	$9.3799 \pm 0.1220$	$2.407 \pm 0.0826$
SWA-2	$0.4100 \pm 0.0135$	$0.6146 \pm 0.0077$	$0.7254 \pm 0.0055$	$0.7480 \pm 0.1373$	$9.5032 \pm 0.1696$	$2.2632 \pm 0.0623$
SWA-3	$0.4869 \pm 0.0070$	$0.6805 \pm 0.0080$	$0.7947 \pm 0.0076$	$0.6248 \pm 0.0621$	$9.6316 \pm 0.1798$	$2.6038 \pm 0.1162$
SWA-4	$0.4057 \pm 0.0063$	$0.6029 \pm 0.0070$	$0.7211 \pm 0.0037$	$0.9196 \pm 0.1833$	$9.3885 \pm 0.1943$	$2.4265 \pm 0.1204$
NA-16	$0.3732 \pm 0.0121$	$0.5723 \pm 0.0093$	$0.6869 \pm 0.0075$	$1.1105 \pm 0.1239$	$9.1334 \pm 0.1387$	$2.6973 \pm 0.0513$
NA-32	$0.3857 \pm 0.0121$	$0.5814 \pm 0.0072$	$0.6973 \pm 0.0071$	$0.6952 \pm 0.0600$	$9.2068 \pm 0.1127$	$2.5957 \pm 0.0464$

SWA = Sliding Window Attention

1. Window size = 4

2. Window size = 8

3. Window size = 12

4. Window size = 16

NA = Neighborhood Attention

	R Precision (top 1)↑	R Precision (top 2)↑	R Precision (top 3)↑	FID ↓	Diversity →	Multimodality ↑
SWA-4	$0.4057 \pm 0.0063$	$0.6029 \pm 0.0070$	$0.7211 \pm 0.0037$	$0.9196 \pm 0.1833$	$9.3885 \pm 0.1943$	$2.4265 \pm 0.1204$
tubelet_p4o0	$0.4121 \pm 0.0158$	$0.6109 \pm 0.0205$	$0.7240 \pm 0.0161$	$0.6275 \pm 0.0840$	$9.5540 \pm 0.1772$	$2.2371 \pm 0.0495$
tubelet_p8o4	$0.4064 \pm 0.0112$	$0.6092 \pm 0.0089$	$0.7209 \pm 0.0054$	$0.8067 \pm 0.0866$	$9.2741 \pm 0.2226$	$2.4800 \pm 0.0705$
convk5s2	$0.4111 \pm 0.0147$	$0.6086 \pm 0.0133$	$0.7232 \pm 0.0079$	$0.9910 \pm 0.1098$	$9.2075 \pm 0.1155$	$2.4582 \pm 0.1058$
convk7s3	$0.4062 \pm 0.0064$	$0.6041 \pm 0.0181$	$0.7166 \pm 0.0137$	$0.9615 \pm 0.1339$	$9.3084 \pm 0.1009$	$2.4867 \pm 0.1116$

Window size = 16