

Method	CIFAR-20	CelebA-HQ	ArtBench (Post-Impressionism)
Model behavior	Inception score	Diversity entropy	Aesthetic score
Pixel similarity (average)	-11.81 ± 4.56	-8.91 ± 0.93	11.24 ± 0.63
Pixel similarity (max)	-31.80 ± 2.90	21.70 ± 2.05	14.61 ± 2.72
Embedding dist. (average)	-	13.83 ± 1.12	-
Embedding dist. (max)	-	7.32 ± 3.16	-
CLIP similarity (average)	5.79 ± 3.67	-32.23 ± 0.87	-6.96 ± 4.08
CLIP similarity (max)	11.31 ± 0.37	-0.93 ± 3.83	-1.75 ± 4.07
Gradient similarity (average)	5.79 ± 3.67	-18.32 ± 0.65	0.25 ± 1.18
Gradient similarity (max)	-0.89 ± 3.17	-12.90 ± 1.60	10.48 ± 3.11
Aesthetic score (average)	-	-	24.85 ± 2.30
Aesthetic score (max)	-	-	21.36 ± 3.70
Relative IF	5.23 ± 5.50	-1.07 ± 0.68	-5.02 ± 1.77
Renormalized IF	11.39 ± 6.79	10.17 ± 0.57	-11.41 ± 0.93
TRAK	7.94 ± 5.67	3.22 ± 0.75	-8.18 ± 1.30
Journey-TRAK	-42.92 ± 2.15	-2.88 ± 4.02	-11.41 ± 4.22
D-TRAK	10.90 ± 1.21	-27.23 ± 2.80	11.30 ± 3.47
Leave-one-out (LOO)	30.66 ± 6.11	-1.22 ± 6.34	3.74 ± 8.00
Sparsified-FT Shapley (Ours)	61.48 ± 2.27	26.34 ± 3.42	61.44 ± 2.04