

Method	CLIP		Human Evaluation		
	Text-Align	Frame-Con	Edit-Acc	Preserve-Acc	Frame-Con
Tune-A-Video (Wu et al. 2022)	0.810	0.959	2.99	3.13	3.05
Control-A-Video (Chen et al. 2022)	0.801	0.955	2.25	2.50	2.88
ControlVideo (Zhang et al. 2022)	0.822	0.963	2.36	2.02	2.08
Gen-1 (Esser et al. 2022)	0.833	0.939	2.41	2.51	2.56
Ground-A-Video (Ours)	0.837	0.970	4.13	4.24	4.06