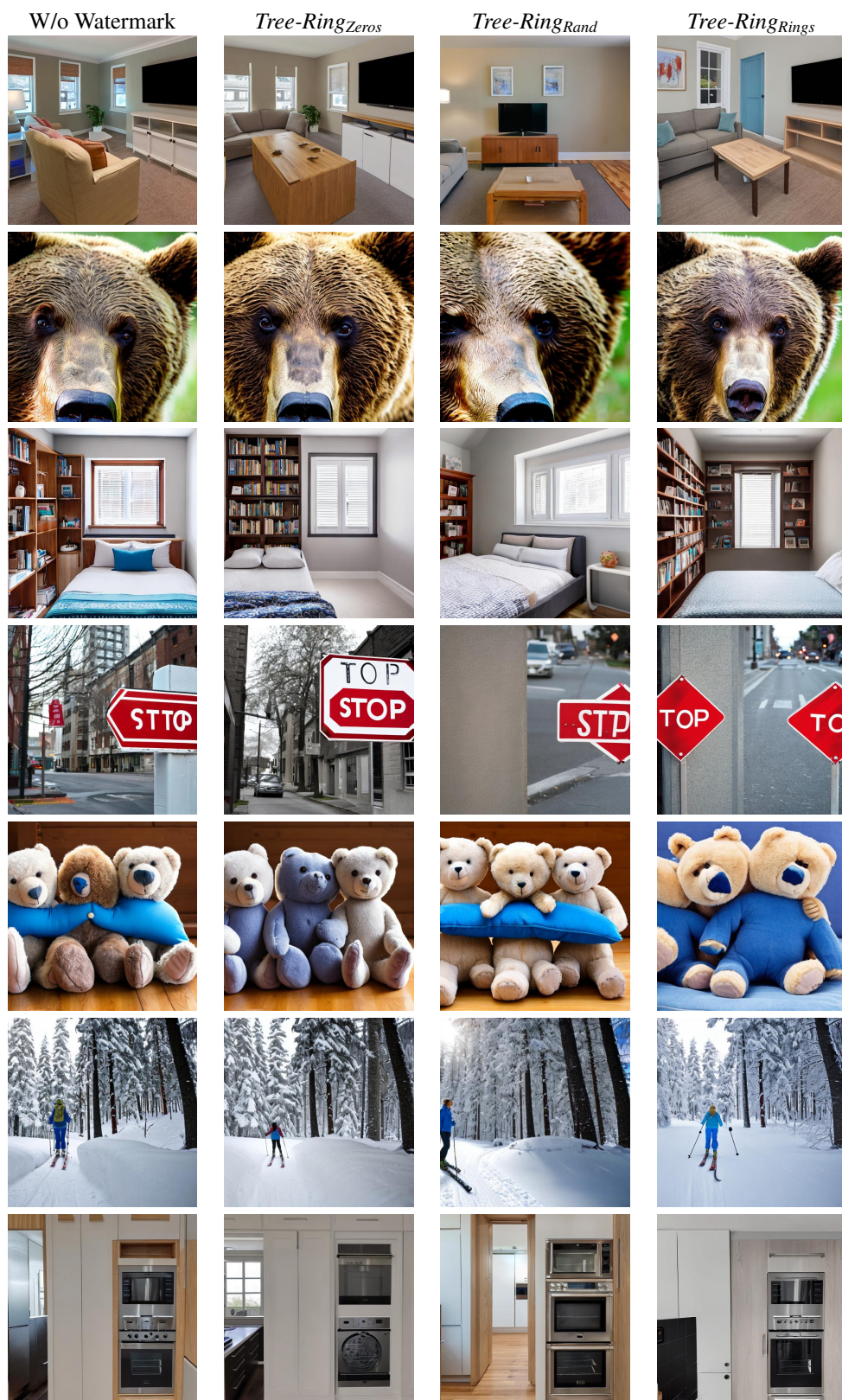


**Table 3:** Main Results with Error Bars. T@1%F represents TPR@1%FPR. We evaluate watermark accuracy in both benign and adversarial settings. Adversarial here refers to average performance over a battery of image manipulations.

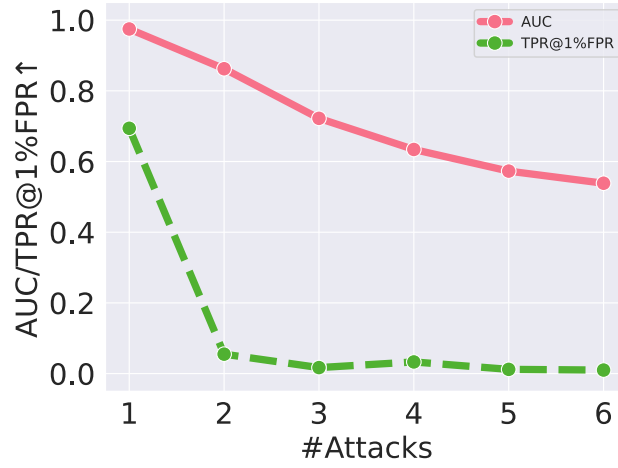
Model	Method	AUC/T@1%F (Clean)	AUC/T@1%F (Adversarial)	FID ↓	CLIP Score ↑
<b>Stable Diff.</b> FID = 25.29 CLIP Score = 0.363	DwtDct	0.974 <sub>.001</sub> / 0.624 <sub>.013</sub>	0.574 <sub>.005</sub> / 0.092 <sub>.004</sub>	25.10 <sub>.09</sub>	0.362 <sub>.000</sub>
	DwtDctSvd	1.000 <sub>.000</sub> / 1.000 <sub>.000</sub>	0.702 <sub>.000</sub> / 0.262 <sub>.011</sub>	25.01 <sub>.09</sub>	0.359 <sub>.000</sub>
	RivaGAN	0.999 <sub>.000</sub> / 0.999 <sub>.000</sub>	0.854 <sub>.002</sub> / 0.448 <sub>.006</sub>	<b>24.51<sub>.17</sub></b>	0.361 <sub>.000</sub>
	<b><i>T-R<sub>Zeros</sub></i></b>	0.999 <sub>.000</sub> / 0.999 <sub>.000</sub>	0.963 <sub>.001</sub> / <b>0.715<sub>.021</sub></b>	26.56 <sub>.07</sub>	0.356 <sub>.000</sub>
	<b><i>T-R<sub>Rand</sub></i></b>	1.000 <sub>.000</sub> / 1.000 <sub>.000</sub>	0.918 <sub>.005</sub> / 0.702 <sub>.017</sub>	25.47 <sub>.05</sub>	0.363 <sub>.001</sub>
	<b><i>T-R<sub>Rings</sub></i></b>	1.000 <sub>.000</sub> / 1.000 <sub>.000</sub>	<b>0.975<sub>.001</sub></b> / 0.694 <sub>.018</sub>	25.93 <sub>.13</sub>	<b>0.364<sub>.000</sub></b>
<b>ImageNet</b> FID = 17.73	DwtDct	0.899 <sub>.040</sub> / 0.244 <sub>.203</sub>	0.536 <sub>.016</sub> / 0.037 <sub>.029</sub>	17.77 <sub>.01</sub>	-
	DwtDctSvd	1.000 <sub>.000</sub> / 1.000 <sub>.000</sub>	0.713 <sub>.019</sub> / 0.187 <sub>.008</sub>	18.55 <sub>.02</sub>	-
	RivaGAN	1.000 <sub>.000</sub> / 1.000 <sub>.000</sub>	0.882 <sub>.010</sub> / 0.509 <sub>.009</sub>	18.70 <sub>.02</sub>	-
	<b><i>T-R<sub>Zeros</sub></i></b>	0.999 <sub>.000</sub> / 1.000 <sub>.000</sub>	0.921 <sub>.000</sub> / 0.476 <sub>.000</sub>	18.78 <sub>.00</sub>	-
	<b><i>T-R<sub>Rand</sub></i></b>	0.999 <sub>.000</sub> / 1.000 <sub>.000</sub>	0.940 <sub>.004</sub> / 0.585 <sub>.006</sub>	18.68 <sub>.09</sub>	-
	<b><i>T-R<sub>Rings</sub></i></b>	0.999 <sub>.000</sub> / 0.999 <sub>.000</sub>	<b>0.966<sub>.005</sub></b> / <b>0.603<sub>.006</sub></b>	<b>17.68<sub>.16</sub></b>	-

**Table 4:** AUC under each Attack for the ImageNet model, showing the effectiveness of *Tree-RingRings* over a number of augmentations. Cr. & Sc. refers to random cropping and rescaling.

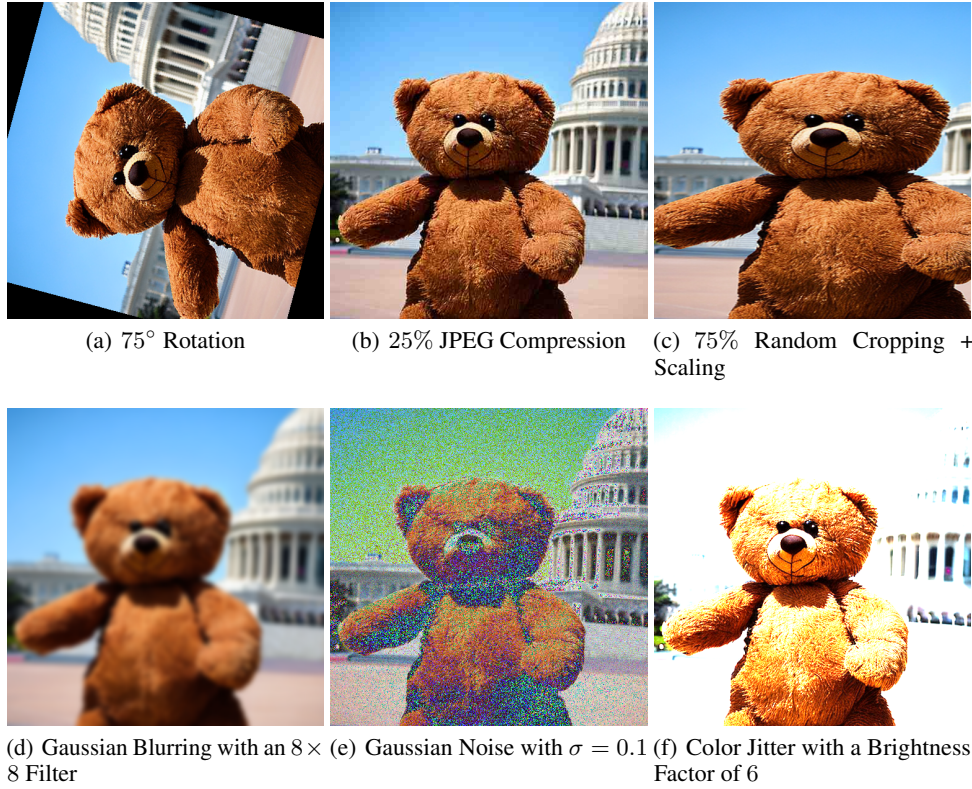
Method	Clean	Rotation	JPEG	Cr. & Sc.	Blurring	Noise	Color Jitter	Avg
DwtDct	0.899	0.478	0.522	0.433	0.512	0.365	0.538	0.536
DwtDctSvd	<b>1.000</b>	0.669	0.568	0.614	0.947	0.656	0.535	0.713
RivaGan	<b>1.000</b>	0.321	<b>0.978</b>	<b>0.999</b>	0.988	0.962	0.924	0.882
<b><i>T-R<sub>Zeros</sub></i></b>	0.999	0.953	0.806	0.997	<b>0.999</b>	0.938	0.775	0.921
<b><i>T-R<sub>Rand</sub></i></b>	0.999	0.682	0.962	0.997	<b>0.999</b>	<b>0.986</b>	<b>0.956</b>	0.940
<b><i>T-R<sub>Rings</sub></i></b>	0.999	<b>0.975</b>	0.940	0.994	<b>0.999</b>	0.979	0.861	<b>0.966</b>



**Figure 6:** More generated images with *Tree-Ring Watermarking* with the first 7 prompts in MS-COCO-2017 training dataset.



**Figure 7:** Results on  $k$  number of random attacks applied at the same time.



**Figure 8:** Attacked images.