








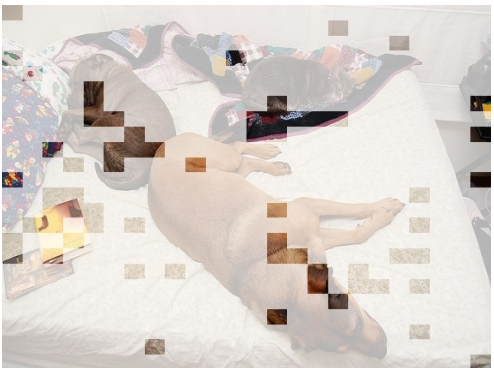






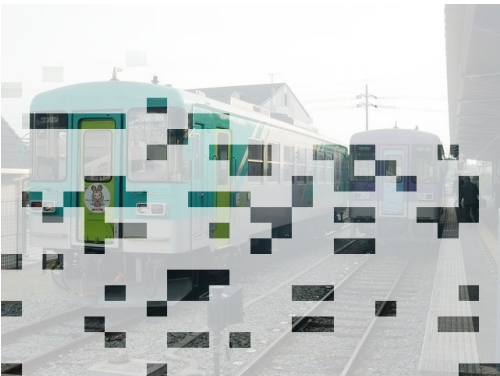
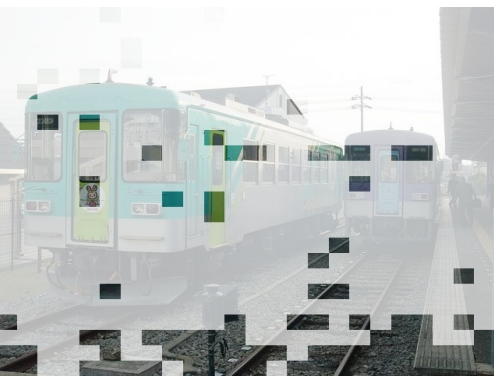

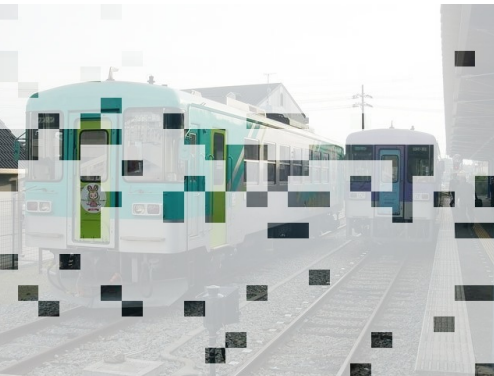

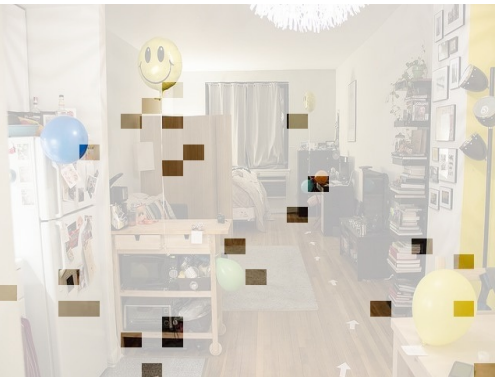
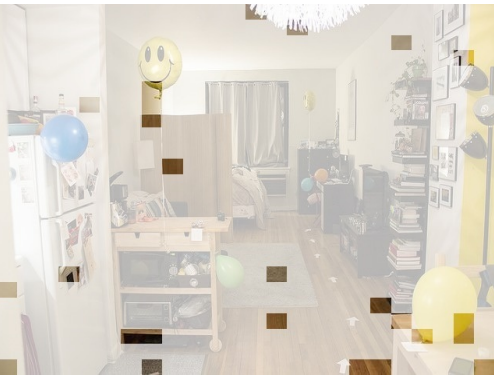
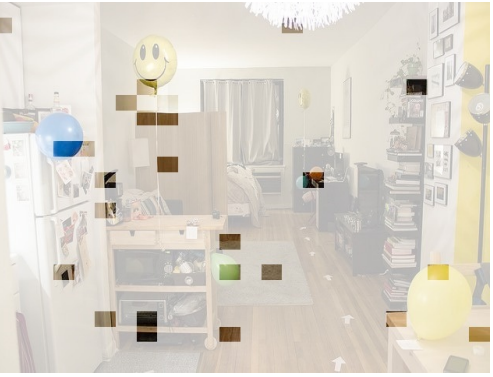
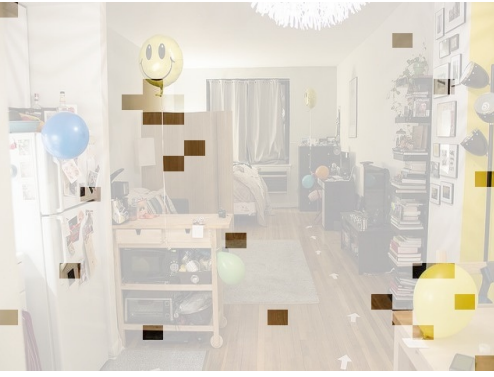


<p>Q:Is the word "DRINKS"placed below "COLD" or above it?</p>		(a)	(b)	(c)	(d)
Input Image		Target-Driven(ours)	Text to Image Attention	Diversity	Cls-token Attention
					
	<p>A:"Below" ✓</p>	<p>A:"Below" ✓</p>	<p>A:"Above" ✗</p>	<p>A:"Above" ✗</p>	<p>A:"Above" ✗</p>
<p>Q:Output a single digit only: how many animals are on the bed?</p>		(a)	(b)	(c)	(d)
Input Image		Target-Driven(ours)	Text to Image Attention	Diversity	Cls-token Attention
					
	<p>A:"3" ✓</p>	<p>A:"3" ✓</p>	<p>A:"2" ✗</p>	<p>A:"2" ✗</p>	<p>A:"2" ✗</p>
<p>Q:Output a single digit only: how many small condiment cups contain dark red jam?</p>		(a)	(b)	(c)	(d)
Input Image		Target-Driven(ours)	Text to Image Attention	Diversity	Cls-token Attention
					
	<p>A:"2" ✓</p>	<p>A:"2" ✓</p>	<p>A:"1" ✗</p>	<p>A:"1" ✗</p>	<p>A:"1" ✗</p>
<p>Q:Answer with exactly one token (YES/ NO): is there a person standing on the platform?</p>		(a)	(b)	(c)	(d)
Input Image		Target-Driven(ours)	Text to Image Attention	Diversity	Cls-token Attention
					
	<p>A:"Yes" ✓</p>	<p>A:"Yes" ✓</p>	<p>A:"No" ✗</p>	<p>A:"No" ✗</p>	<p>A:"No" ✗</p>
<p>Q:Output a single digit only: how many balloons are tied to the stand near the back desk?</p>		(a)	(b)	(c)	(d)
Input Image		Target-Driven(ours)	Text to Image Attention	Diversity	Cls-token Attention
					
	<p>A:"3" ✓</p>	<p>A:"3" ✓</p>	<p>A:"0" ✗</p>	<p>A:"2" ✗</p>	<p>A:"2" ✗</p>