

Self-Corrected Image Generation with Explainable Latent Rewards

Supplementary Material

001	Contents			
002	A Method Details	1		
003	A.1. Transformer-Based Corrector Architecture	1		
004	A.2. Transformer-Based Latent Reward Projection	1		
005	B Additional Qualitative Results	2		
006	B.1. Text-to-Image Generation	2		
007	B.2. Editing Tasks	2		
008	C Supervised Data Generation	2		
009	D Implementation and Reproducibility	2		
010	D.1. Training Configuration	2		
011	D.2. Hyperparameters	2		
012	E Additional Experiments	2		
013	E.1. Process-Level Validation and Attribution	2		
014	E.2. Sub-Reward Analysis and Ablation	6		
015	E.3. Comparison with Plug-and-Play Guidance			
016	Baseline	6		
017	E.4. Confidence-Based Modulation (CMD) Analysis	6		
018	E.5. Additional Technical Details	6		
019	E.6. Runtime and Memory	7		
020				
021	A. Method Details			
022	A.1. Transformer-Based Corrector Architecture			
023	The Understanding-Guided Reinforcement Corrector			
024	(URC) is implemented as a lightweight transformer that			
025	operates on the latent feature grid of the frozen text-to-			
026	image model. Given a latent representation $z_0 \in \mathbb{R}^{C \times H \times W}$			
027	and the prompt embedding e_p , URC predicts a residual			
028	$\Delta_\theta(z_0, e_p)$ applied before decoding.			
029	Latent Tokenization. We flatten the spatial dimensions			
030	and treat each spatial location as a token: $Z_0 =$			
031	$\text{reshape}(z_0) \in \mathbb{R}^{(HW) \times C}$. 2D sine-cosine positional en-			
032	codings are added to retain spatial structure.			
033	Prompt Conditioning. The prompt embedding $e_p \in \mathbb{R}^d$			
034	is projected to the latent dimension via $e'_p = W_p e_p$			
035	and integrated using Feature-wise Linear Modulation (FiLM):			
036	$Z'_0 = \gamma(e'_p) \odot Z_0 + \beta(e'_p),$			
037	where γ and β are produced from e'_p through small MLPs.			
038	This allows the transformer to modulate latent tokens ac-			
039	cording to the semantic content of the prompt.			
	Transformer Layers. URC consists of six transformer			040
	blocks, each containing multi-head self-attention with 8			041
	heads, cross-attention to prompt tokens, feedforward net-			042
	works (hidden size $4C$), and pre-norm residual connec-			043
	tions. Cross-attention explicitly injects linguistic structure,			044
	including object categories, colors, and relational phrases.			045
	Residual Prediction. The output of the transformer is			046
	projected via a linear layer W_o back to latent dimension-			047
	ality, reshaped to form the residual:			048
	$\Delta_\theta(z_0, e_p) \in \mathbb{R}^{C \times H \times W}.$			049
	Despite its transformer architecture, URC remains compact			050
	($\leq 15M$ parameters), ensuring it refines semantics without			051
	overpowering the frozen generator.			052
	A.2. Transformer-Based Latent Reward Projection			053
	The latent reward projector R_ϕ is a transformer that maps			054
	corrected latent activations and prompt embeddings to in-			055
	terpretable rewards approximating image-level feedback.			056
	Input Construction. R_ϕ receives the corrected latent to-			057
	kens $Z_c \in \mathbb{R}^{(HW) \times C}$, the prompt token embeddings $\{e_{p,i}\}$,			058
	and the CMD-derived global semantic vector $g_{\text{cmd}} \in \mathbb{R}^d$,			059
	which is appended as an extra token:			060
	$X_0 = [Z_c; e_{p,1}; \dots; e_{p,T}; g_{\text{cmd}}].$			061
	Transformer Design. R_ϕ has four transformer layers			062
	with 8-head multi-head attention, alternating self-attention			063
	and cross-attention, feedforward networks of size $4C$, and			064
	rotary positional embeddings for prompt tokens.			065
	Reward Heads. After the transformer, the updated se-			066
	matic token g'_{cmd} is passed through three linear layers to			067
	produce the latent reward vector:			068
	$r_{\text{latent}} = [W_{\text{count}} g'_{\text{cmd}}, W_{\text{color}} g'_{\text{cmd}}, W_{\text{pos}} g'_{\text{cmd}}] \in \mathbb{R}^3,$			069
	corresponding to counting, color, and position sub-rewards.			070
	Training Objective. The projector is trained to mini-			071
	mize the L2 distance between predicted latent rewards and			072
	image-level task-specific rewards:			073
	$\mathcal{L}_{\text{proj}} = \sum_{i=1}^3 \ r_{\text{latent}}^{(i)} - r_{\text{image}}^{(i)}\ _2^2,$			074
	enabling gradient-based optimization of URC even when			075
	the original image-level reward is non-differentiable.			076

077	B. Additional Qualitative Results		
078	B.1. Text-to-Image Generation		
079	Additional qualitative results for text-to-image generation		
080	using our model are shown in Fig.1. The prompts used for		
081	text-to-image generation, starting from top-left and going		
082	row-wise from left to right, are as follows:		
083	1. A person walking alone on a quiet street at sunset.		
084	2. A bowl of fresh fruit sitting on a kitchen counter.		
085	3. A dog lying on a couch in a cozy living room.		
086	4. A car parked beside a forest road in the morning.		
087	5. A cup of coffee on a wooden table near a window.		
088	6. A small boat floating on a calm lake at dawn.		
089	7. A cyclist riding through a city park.		
090	8. A marketplace stall filled with colorful vegetables.		
091	9. A cat sitting on a windowsill looking outside.		
092	10. A person reading a book in a quiet café.		
093	11. A train passing through a snowy landscape.		
094	12. A street food vendor cooking at night.		
095	B.2. Editing Tasks		
096	The prompts used to get results for the image editing task		
097	Figure 2 Figure 3 Figure 4, starting from the top-left and		
098	going row-wise from left to right, are as follows:		
099	1. Add a small backpack resting on the ground next to the		
100	bicycle.		
101	2. Add a small cushion under the cat.		
102	3. Turn the baker into a man.		
103	4. Add a small cup next to the pitcher.		
104	5. Add a small dog walking next to the couple.		
105	6. Remove half the people from the image, from the cross-		
106	walk.		
107	7. Change the kite's color to bright red.		
108	8. Place a small cooler next to the fisherman, near his feet.		
109	9. Add a line of people near the food truck.		
110	10. Add a straw to the glass.		
111	11. Add an umbrella above the table.		
112	12. Add a plate of vegetables on the grill.		
113	13. Add a notebook next to the laptop.		
114	14. Add a small trail sign beside the path.		
115	15. Add a small seashell next to the sneakers.		
116	16. Add a distant mountain in the background scenery.		
117	17. Add a small bench next to the bus stop.		
118	18. Add a single cloud in the sky above the mountain.		
119	19. Add a small water bottle on the ground next to the per-		
120	son.		
121	20. Add a few blueberries on the plate beside the pancakes.		
122	21. Add a folded blanket at the side of the bed.		
123	22. Add an open guitar case on the ground in front of him.		
124	23. Add a price tag to one of the flower pots.		
125	24. Remove the plant pots from the row.		
	C. Supervised Data Generation		126
	Training prompts were generated using large language		127
	models to cover three semantic domains: color, position,		128
	and object count. For each domain, 10k candidate prompts		129
	were generated using multiple models (BLIP3o, SD3, and		130
	DEV). To ensure high-quality and semantically accurate		131
	prompts, all outputs were manually reviewed and filtered,		132
	retaining only those that correctly captured the intended		133
	attributes. Among the models tested, BLIP3o, SD3, and		134
	DEV consistently produced the most reliable and coherent		135
	prompts, showing strong consistency in describing colors,		136
	spatial relations, and object counts, whereas other models		137
	occasionally produced ambiguous or incomplete descrip-		138
	tions.		139
	Rather than using a separate testing set, the supervised		140
	dataset is used to guide the generator itself: for each val-		141
	idated prompt, the model generates multiple images using		142
	different random seeds. This self-generation process, super-		143
	vised by the validated dataset, provides diverse latent rep-		144
	resentations and ensures coverage of the semantic domains,		145
	without requiring an explicit test split. The resulting dataset		146
	thus serves both as a source of training supervision and a		147
	reference for controlled evaluation during model develop-		148
	ment.		149
	D. Implementation and Reproducibility		150
	D.1. Training Configuration		151
	We train using AdamW with a learning rate of 1×10^{-4} ,		152
	batch size 8 per GPU, PPO clipping ratio 0.2, gradient clip-		153
	ping 1.0, and a cosine learning rate schedule. Experiments		154
	are conducted on H100 80GB GPUs.		155
	D.2. Hyperparameters		156
	The URC transformer has 6 layers, R_ϕ has 4 layers, em-		157
	bedding size 1024, residual scale $\alpha = 0.8$, and task-weight		158
	modulation range $[0.5, 2.0]$.		159
	E. Additional Experiments		160
	This section provides additional quantitative experiments		161
	complementing the analyses in the main paper.		162
	E.1. Process-Level Validation and Attribution		163
	Causal Evidence. In Sec. 4.3.1, we demonstrate that im-		164
	provements are directly driven by our formulation via three		165
	validation tests: (i) Spatial Causal Link: Masking high-		166
	activation LAM regions leads to a 6.3% drop in CLIP-		167
	Score, proving these regions are where alignment is fixed;		168
	(ii) Signal-to-Gain Correlation: A Spearman correlation		169
	of $\rho=0.71$ between token contribution magnitudes and re-		170
	ward gains confirms that latent rewards directly drive the		171
	correction; (iii) Consistency: A Jaccard similarity of 0.68		172



Figure 1. Qualitative results for text to image generation.

173 across prompts shows stable, predictable correction pat-
174 terns.

175 **Incremental Dynamics.** To quantify process-level dynam-
176 ics, we analyzed semantic trajectories on a subset of 500

prompts from GenEval by sampling latent sub-rewards at 177
irregular intervals through our differentiable projector R_{Φ} . 178
We observe a 26.34% average reward increase between de- 179
noising steps 7 and 38, with a high Pearson correlation 180

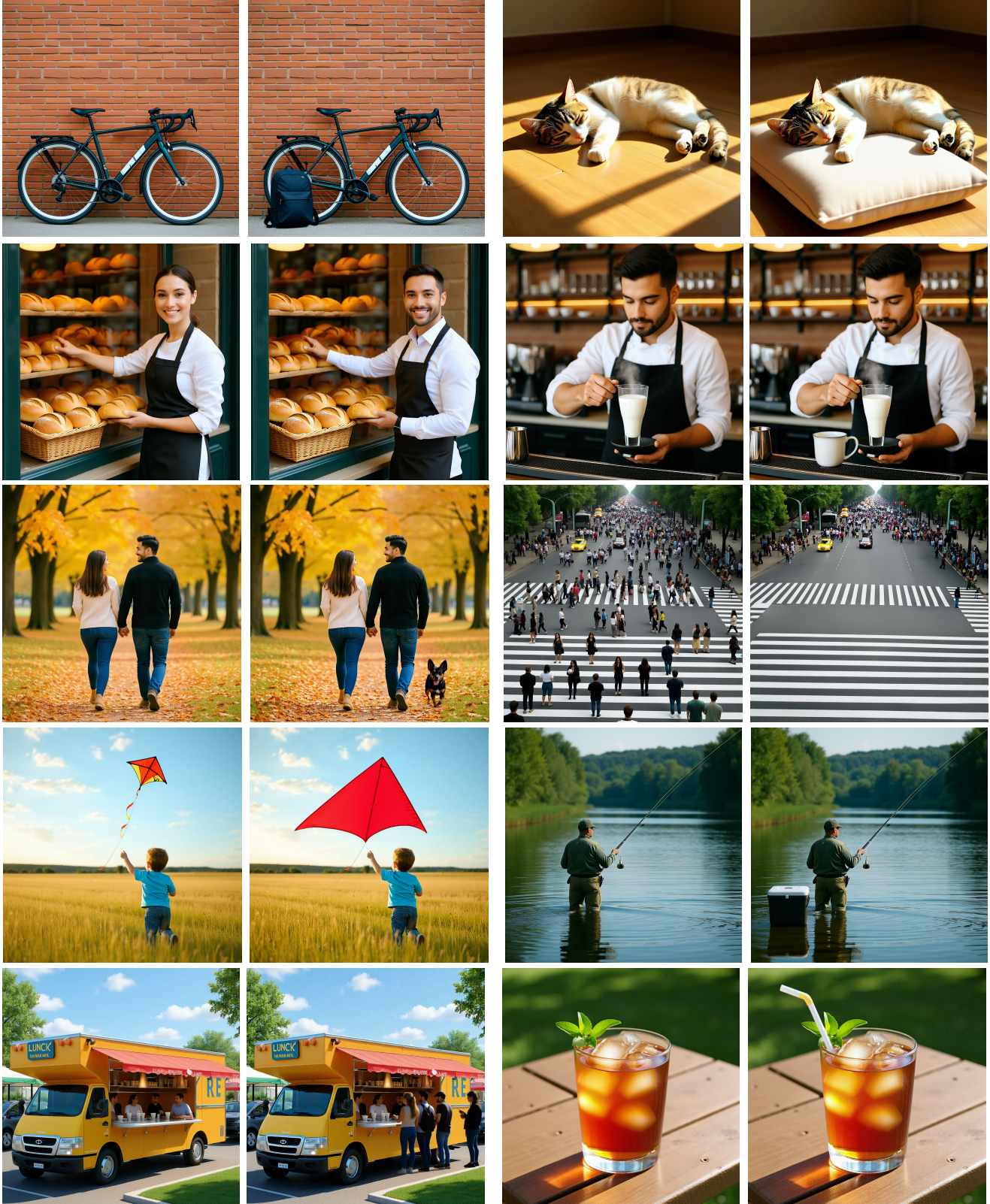


Figure 2. Qualitative editing results-1.

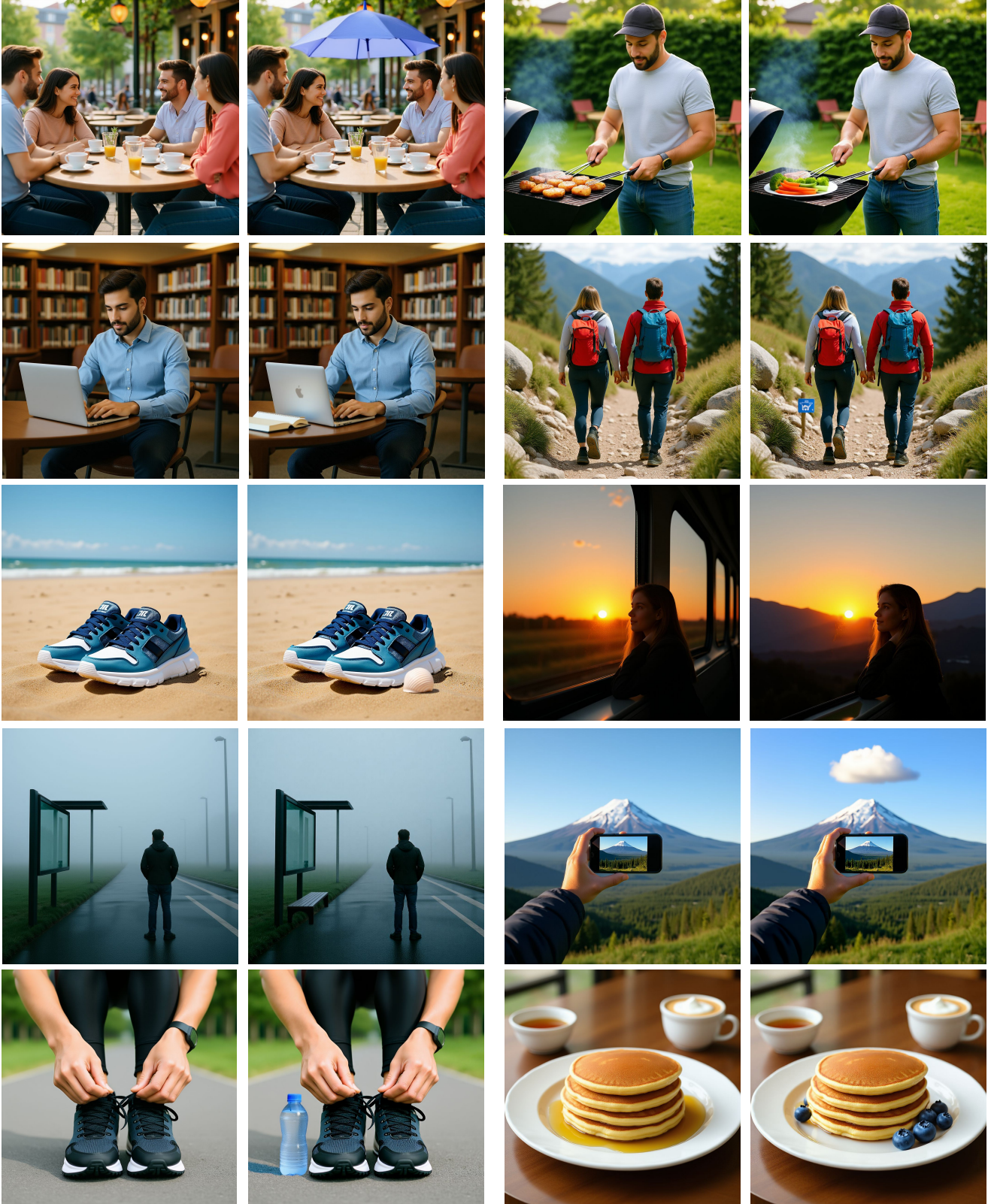


Figure 3. Qualitative editing results-2.



Figure 4. Qualitative editing results-3.

Table 1. Ablation study of reward components.

Method Variant	GenEval (%)	DPG-Bench (%)
Full xLARD Framework	81.24	86.72
w/o Counting Reward	76.45 (−4.79)	83.10 (−3.62)
w/o Color Reward	78.92 (−2.32)	84.55 (−2.17)
w/o Position Reward	79.15 (−2.09)	84.88 (−1.84)

($r=0.827$) between these incremental gains and final performance. Specifically, Counting rewards stabilize rapidly within the first 18 steps to establish structural layout, while Color and Position rewards provide continuous refinement through step 43. This temporal analysis confirms that xLARD performs active, multi-stage steering throughout the generative process rather than a singular post-hoc correction.

E.2. Sub-Reward Analysis and Ablation

Our modular design is a deliberate choice to ensure interpretability and avoid the black-box nature of aggregate rewards. While the sub-rewards are specific, they address the most frequent failure modes in T2I models and provide a framework that is easily extendable to other attributes. The ablation of reward roles is reported in Table 1.

E.3. Comparison with Plug-and-Play Guidance Baseline

For latent editing comparisons, we include a Plug-and-Play Guidance (PPGD) baseline adapted to our architecture; most existing latent editing methods are U-Net-based

and do not generalize to unified multimodal transformers. On GenEval and DPG-Bench, PPGD achieves 77.04% and 83.54%, below our method (81.29% and 86.45%). These results provide additional comparative context against recent SOTA techniques.

E.4. Confidence-Based Modulation (CMD) Analysis

The Confidence Head ω is a lightweight MLP trained via PPO to predict sub-reward reliability. Rather than static weights, it acts as a dynamic gating mechanism to suppress irrelevant signals. Our analysis shows removing CMD causes a 3.3% GenEval drop due to “gradient interference,” where unmodulated rewards compete for latent updates.

E.5. Additional Technical Details

Reward Projection. The differentiable projector R_ϕ does not backpropagate through the decoder. Instead, it is trained via supervised regression to approximate non-differentiable image-level rewards, enabling stable gradient flow entirely within latent space during corrector optimization.

Corrector Behavior. The corrector Δ_θ is explicitly constrained to produce small-magnitude residual updates via the scaling factor α and PPO regularization, ensuring localized semantic refinement rather than latent overwriting or re-generation.

Counting Implementation. Connected-component analysis is applied after adaptive thresholding and morphological filtering of token attention maps, which removes spurious

228 activations and yields stable object-count estimates across
229 prompts.

230 **E.6. Runtime and Memory**

231 Training takes 4 h/epoch, inference matches base-generator
232 runtime as mentioned in the paper discussion, and peak
233 memory is approximately 72 GB at batch size 8. Reference-
234 prompt construction details are provided in the supplemen-
235 tary material.