

# Towards Fine-Grained Attribution: Instance-Aware Preference Optimization for Aligning Diffusion Models

## Supplementary Material

### A. Additional Implementation details

To ensure the fairness and reproducibility of all experiments, this section provides additional key implementation details. Regarding model initialization, for Stable Diffusion 1.5 (SD1.5) [7], both IAPO and InPO [3] utilize the same high-quality training starting point: the officially provided supervised fine-tuned (SFT) model weights available in their Hugging Face repository. All other baseline models (DPO [6], KTO [2]) directly employ the official pre-trained checkpoints released by their respective authors on Hugging Face, ensuring all comparisons are conducted on an authoritative and consistent baseline. For Stable Diffusion XL (SDXL) [4], to enable the most rigorous comparison and highlight the advantage of our proposed instance-level dataset, we fully reproduce the official InPO codebase and train both IAPO and InPO using identical hyperparameters. During evaluation, the inference configuration remains unified across all text-to-image generation: the classifier-free guidance scale is set to 7.5 for SD1.5 and 5.0 for SDXL, with sampling steps fixed at 50. The training random seed is fixed at 42 for all training models, while the test-time random seed is fixed at 0. We provide the complete evaluation results for SDXL with seed 66 Tab. I to demonstrate the generalization capability of IAPO. Notably, the performance variance of IAPO remains minimal, demonstrating exceptional training stability. This consistent margin across multiple seeds confirms that our method’s advantages are intrinsic to its design rather than artifacts of stochastic variation, thereby reinforcing the reliability of our main conclusions.

### B. Training Efficiency Analysis on SDXL

In Fig. I, we present comparative results of IAPO and baseline methods in terms of image quality and training efficiency based on SDXL. As shown in Fig. I, IAPO reaches training speeds that are  $6.25\times$  faster than Diffusion-DPO, and  $1.31\times$  faster than InPO, while also generating higher-quality images.

### C. Introduction to the Dataset

**Pick-a-Pic v2** dataset is a large-scale collection of text-to-image pairs annotated with human preferences, sourced from the Pick-a-Pic web application. Each data point is structured as a tuple,  $(caption, jpg_0, jpg_1, label)$ , containing a text prompt, two corresponding generated images,

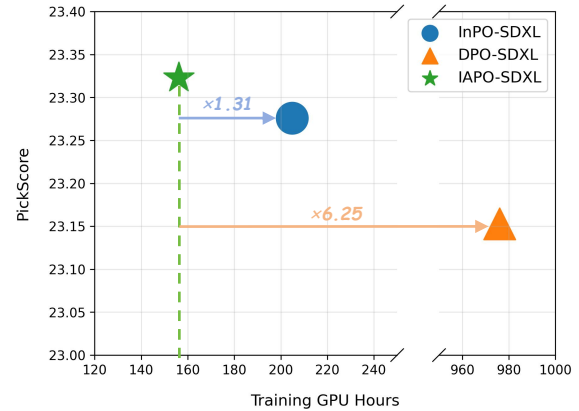


Figure I. A comparative evaluation of image quality and training efficiency between IAPO and baselines of SDXL on HPD v2.

and a binary label indicating the user’s choice. The visual corpus was synthesized using a variety of text-to-image models—including Stable Diffusion 2.1, Dreamlike Photoreal, and several variants of Stable Diffusion XL (SDXL)—across a wide spectrum of Classifier-Free Guidance (CFG) scales.

**HPD v2** dataset is curated by collecting human preferences on the "Dreambot" channel within the Stable Foundation Discord server. It contains 98,807 images generated from 25,205 prompts. The data structure links each prompt to a variable number of images, which are organized into pairs with labels denoting the user’s preferred choice. For our analysis, we use the 3,200 prompts from the official test split.

**Parti-Prompts** dataset includes 1,632 text prompts created to evaluate text-to-image generation models. These prompts are divided into multiple categories to provide a wide range of challenges, which helps to assess model performance thoroughly across different areas.

**DrawBench** serves as a benchmark for assessing text-to-image generation models across 11 specific categories, such as object composition, spatial relationships, style consistency, and text accuracy. Through carefully crafted prompts, it evaluates a model’s ability to translate semantic and stylistic cues into visuals, focusing on fidelity, creativity, and alignment with descriptions. By using side-by-side comparisons, DrawBench offers a qualitative approach to identifying strengths and weaknesses, fostering advancements in AI-driven generative models and computer vision.

Table I. We assess the stability of IAPO under multiple random seeds by fine-tuning SDXL and evaluating on Parti-Prompts, HPD v2 and Pick-a-Pic v2. Results are reported as mean and median scores; the best value for each metric is **bolded**.

Datasets	Models	Aesthetic [8]		PickScore [1]		HPS [9]		CLIP [5]	
		Mean	Median	Mean	Median	Mean	Median	Mean	Median
Parti-Prompts [10]	SDXL	5.7693	5.7571	22.6274	22.6496	28.4241	28.4723	35.5172	35.5675
	InPO (Seed=42)	5.8332	5.8150	22.9952	<b>22.9889</b>	28.8443	<b>28.8621</b>	<b>35.6453</b>	<b>35.6660</b>
	IAPO (Seed=42)	<b>5.9503</b>	<b>5.9439</b>	<b>23.0246</b>	22.9863	<b>28.8745</b>	28.8230	35.3439	35.3803
	InPO (Seed=66)	5.7592	5.7784	22.9152	22.9061	28.8096	28.8425	<b>35.8458</b>	<b>36.0351</b>
	IAPO (Seed=66)	<b>5.8767</b>	<b>5.8740</b>	<b>22.9360</b>	<b>22.9367</b>	<b>28.9871</b>	<b>28.9007</b>	35.4458	35.4845
HPD v2 [9]	SDXL	6.1338	6.1192	22.7835	22.7657	28.6278	28.6167	38.1623	38.3686
	InPO (Seed=42)	6.1797	6.1666	23.2761	23.2338	29.2413	29.2557	38.2722	38.5370
	IAPO (Seed=42)	<b>6.2580</b>	<b>6.2622</b>	<b>23.3230</b>	<b>23.2793</b>	<b>29.2848</b>	<b>29.2988</b>	<b>38.3220</b>	<b>38.5907</b>
	InPO (Seed=66)	6.1003	6.0803	23.1599	23.1429	29.1276	29.2130	<b>38.4731</b>	<b>38.6416</b>
	IAPO (Seed=66)	<b>6.1903</b>	<b>6.1789</b>	<b>23.2344</b>	<b>23.2543</b>	<b>29.2665</b>	<b>29.3127</b>	38.1847	38.3365
Pick-a-Pic v2 [1]	SDXL	6.0040	5.9788	22.1655	22.2426	27.9776	28.0080	36.1164	36.5889
	InPO (Seed=42)	6.0416	6.0167	22.5820	22.5573	28.5171	28.6445	36.5457	<b>36.7526</b>
	IAPO (Seed=42)	<b>6.1226</b>	<b>6.1282</b>	<b>22.6806</b>	<b>22.7006</b>	<b>28.6004</b>	<b>28.6707</b>	<b>36.7268</b>	36.7268
	InPO (Seed=66)	6.0551	6.0551	22.6064	22.6037	28.5605	28.2803	<b>37.9502</b>	39.0377
	IAPO (Seed=66)	<b>6.1918</b>	<b>6.2661</b>	<b>22.8349</b>	<b>23.0578</b>	<b>28.8182</b>	<b>28.7883</b>	37.3843	<b>37.3160</b>

## 074 D. Additional Quantitative Results

075 We further conduct extended quantitative studies to complement the results already reported. All baseline models are evaluated with the official released weights, without any re-training or fine-tuning. For reproducibility, the training random seed is fixed at 42 during generation; and the random seed for test is fixed at 0.

081 SDXL experiments are summarized in Table II. We compare our approach with two recent preference-alignment methods, Diffusion-DPO and MAPO. Evaluations are performed on Pick-a-Pic v2, HPD v2 and Parti-Prompts. Our model obtains the highest average score, outperforming all the baselines.

087 Table III reports DrawBench results for both SD 1.5 and SDXL, using the same baselines and metrics as above. Our method achieves the best overall accuracy on both model scales. Results substantiate that IAPO’s improvements extend to scenarios demanding rigorous text-image correspondence, validating its effectiveness for applications where semantic accuracy is critical.

## 094 E. Additional Ablation Studies

095 We present additional ablation studies to dissect the contributions of key components in our Instance-Aware Preference Optimization (IAPO) framework. All experiments are conducted on SD1.5 using the Pick-a-Pic v2 test set unless otherwise specified.

### E.1. Ablation on $w_{pos}$ and $w_{neg}$

The weight assignment for  $w_{pos}$  follows the same establishment method as  $w_{neg}$ . When the preference for the  $n$ -th instance aligns with the global label—meaning we prefer the instance in  $\mathbf{x}_0^w$  (denoted by  $\rho_n = 0$ )—we assign a weight of  $w_{pos}$  to the corresponding bounding boxes to reinforce the learning of these consistent instances:

$$M_n^*(i, j) = \begin{cases} w_{pos} & \text{if } (i, j) \in b_n^* \text{ and } \rho_n = 0 \\ w_{neg} & \text{if } (i, j) \in b_n^* \text{ and } \rho_n = 1 \\ 1 & \text{otherwise} \end{cases} \quad (1)$$

As shown in Table IV, we vary  $w_{neg}$  while fixing  $w_{pos} = 1$ . Performance consistently improves as  $w_{neg}$  decreases, with the best results achieved at  $w_{neg}=0$ , i.e., completely suppressing gradients from conflicting instances in the losing sample. This confirms that ignoring these distractor regions is more effective than down-weighting them, as it prevents the model from reinforcing undesired artifacts. Notably, adjusting  $w_{pos}$  yields diminishing returns, suggesting that penalizing bad regions in positive samples is more critical than amplifying good regions in negative ones. We therefore adopt  $w_{neg} = 0$  and  $w_{pos} = 1$  as our default configuration.

We further set  $w_{neg}$  to -1 to investigate the model’s behavior. Experimental results show that the training process becomes highly unstable under this setting. We argue that this instability stems from an inherent conflict between pixel-level weight assignment and the optimization objective: when the loss weights for certain pixels are nega-

Table II. We compare different preference optimization methods by fine-tuning SDXL and evaluating on Parti-Prompts, HPD v2 and Pick-a-Pic v2. Results are reported as mean and median scores; the best value for each metric is **bolded** and the second-best is underlined.

Datasets	Models	Aesthetic [8]		PickScore [1]		HPS [9]		CLIP [5]	
		Mean	Median	Mean	Median	Mean	Median	Mean	Median
Parti-Prompts	SDXL	5.7693	5.7571	22.6274	22.6496	28.4241	28.4723	<u>35.5172</u>	<u>35.5675</u>
	DPO	5.7946	5.8097	<u>22.9307</u>	<u>22.9264</u>	<b>28.9085</b>	<b>28.8873</b>	<b>36.4665</b>	<b>36.4017</b>
	MaPO	<u>5.8979</u>	<u>5.9017</u>	22.5946	22.5412	28.5551	28.4649	35.2734	35.3595
	IAPO	<b>5.9503</b>	<b>5.9439</b>	<b>23.0246</b>	<b>22.9863</b>	<u>28.8745</u>	<u>28.8230</u>	35.3439	35.3803
HPD v2	SDXL	6.1338	6.1192	22.7835	22.7657	28.6278	28.6167	38.1623	38.3686
	DPO	6.1124	6.1310	<u>23.1330</u>	<u>23.1520</u>	<u>29.1650</u>	<u>29.1740</u>	<b>38.8650</b>	<b>38.7110</b>
	MaPO	<u>6.2416</u>	<u>6.2453</u>	22.8488	22.8169	28.9939	29.0085	38.1587	38.5515
	IAPO	<b>6.2580</b>	<b>6.2622</b>	<b>23.3230</b>	<b>23.2793</b>	<b>29.2848</b>	<b>29.2988</b>	<u>38.3220</u>	<u>38.5907</u>
Pick-a-Pic v2	SDXL	6.0040	5.9788	22.1655	22.2426	27.9776	28.0080	36.1164	36.5889
	DPO	6.0127	6.0168	<u>22.6397</u>	<u>22.6137</u>	<u>28.5933</u>	28.5017	<u>37.3825</u>	<b>37.4003</b>
	MaPO	<b>6.2037</b>	<b>6.2294</b>	22.3179	22.3404	28.5381	<u>28.6262</u>	<b>37.4294</b>	<u>37.3952</u>
	IAPO	<u>6.1226</u>	<u>6.1282</u>	<b>22.6806</b>	<b>22.7006</b>	<b>28.6004</b>	<b>28.6707</b>	36.7268	36.7268

Table III. We compare different preference optimization methods by fine-tuning SD1.5 and SDXL and evaluating on DrawBench. Results are reported as mean and median scores; the best value for each metric is **bolded** and the second-best is underlined.

Base Model	Models	Aesthetic [8]		PickScore [1]		HPS [9]		CLIP [5]	
		Mean	Median	Mean	Median	Mean	Median	Mean	Median
SD1.5	SD1.5	5.2020	5.2677	21.3496	21.4011	27.4474	27.4640	32.5676	33.0921
	DPO	5.3099	5.3774	21.5932	21.7097	27.7069	27.8015	33.2266	33.7624
	KTO	5.3819	5.4458	21.6549	21.8707	28.2696	28.4872	33.4800	34.2055
	InPO	<u>5.3927</u>	<u>5.4927</u>	<u>21.7661</u>	<u>21.8812</u>	<u>28.4733</u>	<u>28.6451</u>	<b>34.2419</b>	<u>34.8507</u>
	IAPO	<b>5.5822</b>	<b>5.6711</b>	<b>21.8486</b>	<b>21.9922</b>	<b>28.5449</b>	<b>28.7431</b>	<u>34.1385</u>	<b>34.8584</b>
SDXL	SDXL	5.5860	5.6435	22.4924	22.6595	28.5996	28.8884	35.3466	<u>36.5069</u>
	DPO	5.6388	<u>5.7231</u>	<u>22.8259</u>	22.9755	<u>29.0795</u>	<u>29.3985</u>	<b>36.6062</b>	<b>37.8796</b>
	MaPO	<u>5.7401</u>	5.7255	22.5230	22.6339	28.8760	29.0263	34.6284	35.2852
	InPO	5.6779	5.6841	22.8186	<u>23.0278</u>	28.9748	29.2258	35.0993	36.0295
	IAPO	<b>5.7658</b>	<b>5.7566</b>	<b>22.9322</b>	<b>23.1392</b>	<b>29.2072</b>	<b>29.4069</b>	<u>36.0807</u>	36.5051

126 tive, the model receives opposing gradient signals within the  
 127 same image. This signal conflict leads to inconsistent pa-  
 128 rameter update directions, disrupts the stability of the train-  
 129 ing process, and ultimately causes convergence failure.

## 130 E.2. Ablation on $\beta$ Coefficient

131 The coefficient  $\beta$  balances preference alignment against  
 132 deviation from the reference model. We test  $\beta \in$   
 133  $\{2000, 3000, 4000, 5000\}$  with all other settings fixed. Ta-  
 134 ble V shows that while IAPO remains relatively stable  
 135 across this range,  $\beta = 2000$  achieves the highest scores  
 136 on PickScore and HPS. A smaller  $\beta$  allows stronger opti-  
 137 mization toward human preferences, which benefits fine-  
 138 grained instance alignment; however, excessively low val-  
 139 ues risk overfitting to noisy annotations. The robustness  
 140 across  $\beta$  values demonstrates that our instance-level credit  
 141 assignment mechanism is the primary driver of performance

gains, rather than precise regularization tuning. Inversion vs.  
 Direct Sampling for SD1.5.

## E.3. Ablation on Inversion Strategy for SD1.5

142 In our main experiments, we note that DDIM inversion pro-  
 143 duces semantically less faithful latents for SD1.5 compared  
 144 to SDXL, leading us to use direct noise sampling for SD1.5  
 145 training. To systematically validate this design choice, we  
 146 conduct an ablation where SD1.5 is trained with 10-step  
 147 DDIM inversion. The results (Table VI) reveal a noticeable  
 148 drop in generation quality and diversity compared to direct  
 149 sampling, confirming that inversion-induced latent distor-  
 150 tion hampers optimization for smaller models.  
 151  
 152  
 153

## F. Additional Qualitative Results

154 Fig. III and Fig. IV displays additional outcomes for addi-  
 155 tional baselines and our IAPO on SD1.5 and SDXL. These  
 156

Table IV. Ablation study on weight assignment for positive and negative regions on Pick-a-Pic v2 dataset. Results are reported as mean and median scores; the best value for each metric is **bolded** and the second-best is underlined.

Dataset	$w_{\text{pos}}$	$w_{\text{neg}}$	Aesthetic [8]		PickScore [1]		HPS [9]		CLIP [5]	
			Mean	Median	Mean	Median	Mean	Median	Mean	Median
Pick-a-Pic v2	1	1	5.7352	5.7802	21.5235	<u>21.6018</u>	27.7578	27.9123	34.6614	<u>35.0903</u>
	2	1	5.7397	5.7660	21.5434	21.5860	27.7716	27.9111	34.6235	34.7461
	1	0.8	5.7384	5.7678	21.5264	21.5968	27.7672	27.9382	34.7159	34.9041
	1	0.5	5.7626	5.8009	21.5294	21.5695	27.7548	27.9478	<u>34.7548</u>	34.7695
	1	0.2	<u>5.7731</u>	<u>5.8126</u>	21.5502	21.5545	27.7894	27.9762	34.6510	34.7344
	2	0	5.7641	5.7865	<u>21.5674</u>	21.5735	<u>27.8429</u>	<u>28.0401</u>	<b>36.9766</b>	<b>36.4228</b>
	1	0	<b>5.7842</b>	<b>5.8168</b>	<b>21.5828</b>	<b>21.6021</b>	<b>27.8492</b>	<b>28.0742</b>	34.6979	<u>35.0903</u>

Table V. Ablation study on the  $\beta$  parameter on Pick-a-Pic v2 dataset. Results are reported as mean and median scores; the best value for each metric is **bolded** and the second-best is underlined.

Dataset	$\beta$	Aesthetic [8]		PickScore [1]		HPS [9]		CLIP [5]	
		Mean	Median	Mean	Median	Mean	Median	Mean	Median
Pick-a-Pic v2	5000	5.7025	5.7333	21.4195	21.4536	27.7656	27.8892	34.3700	34.8469
	4000	5.7199	5.7559	21.4356	21.4611	27.7901	27.9196	34.4997	34.8425
	3000	<u>5.7497</u>	<u>5.7724</u>	<u>21.5032</u>	<u>21.5245</u>	<u>27.8175</u>	<u>28.0183</u>	<u>34.6176</u>	<u>34.9643</u>
	2000	<b>5.7842</b>	<b>5.8168</b>	<b>21.5828</b>	<b>21.6021</b>	<b>27.8492</b>	<b>28.0742</b>	<b>34.6979</b>	<b>35.0903</b>

157 figures offer a comprehensive visual comparison that elu- 157  
 158 cidates the performance disparities between state-of-the-art 158  
 159 models and IAPO. Our method shows significant improve- 159  
 160 ments over baseline models across multiple dimensions in- 160  
 161 cluding aesthetic quality and text-image alignment, with 161  
 162 particularly notable progress in instance-level generation 162  
 163 performance. 163

## 164 G. Limitations

165 Our approach is constrained by the quality of the origi- 165  
 166 nal dataset, which contains a non-negligible proportion of 166  
 167 low-quality image-text pairs, and by the Vision-Language 167  
 168 Model’s tendency to hallucinate, producing occasional mis- 168  
 169 alignments between visual content and textual descrip- 169  
 170 tions; furthermore, our preference-learning framework is 170  
 171 restricted to the instance level and may be too coarse for 171  
 172 tasks that require fine-grained understanding such as distin- 172  
 173 guishing specific parts (e.g., an eagle’s talons) or pixel-level 173  
 174 attributes. The misuse of advanced image generation tech- 174  
 175 nologies presents ongoing threats to societal trust and infor- 175  
 176 mation ecosystems. We consistently advocate for compre- 176  
 177 hensive ethical governance of IAPO throughout its lifecy- 177  
 178 cle to ensure its development remains aligned with digital 178  
 179 ethics standards. 179

## 180 H. Future work

181 We plan to construct a new high-quality dataset by lever- 181  
 182 aging state-of-the-art generative models to produce visually 182

consistent and semantically accurate image-text pairs. To 183  
 mitigate VLM hallucination, we will incorporate a more ro- 184  
 bust annotation pipeline that includes self-reflection mech- 185  
 anisms and cross-verification across multiple models. Fur- 186  
 thermore, we aim to extend our preference modeling from 187  
 the instance level to a more granular hierarchy, enabling 188  
 part-level and even pixel-level preference learning. These 189  
 improvements are expected to enhance both the reliability 190  
 and expressiveness of the learned representations. 191

## 192 References

- 193 [1] Yuval Kirstain, Adam Polyak, Uriel Singer, Shahbuland Ma- 193  
 194 tiana, Joe Penna, and Omer Levy. Pick-a-Pic: An Open 194  
 195 Dataset of User Preferences for Text-to-Image Generation. 195  
 In *NeurIPS*, 2023. 2, 3, 4, 5 196
- 197 [2] Shufan Li, Konstantinos Kallidromitis, Akash Gokul, 197  
 Yusuke Kato, and Kazuki Kozuka. Aligning Diffusion Mod- 198  
 els by Optimizing Human Utility. In *NeurIPS*, 2024. 1 199
- 200 [3] Yunhong Lu, Qichao Wang, Hengyuan Cao, Xierui Wang, 200  
 Xiaoyin Xu, and Min Zhang. InPO: Inversion Preference 201  
 Optimization with Reparametrized DDIM for Efficient Dif- 202  
 fusion Model Alignment. In *CVPR*, 2025. 1 203
- 204 [4] Dustin Podell, Zion English, Kyle Lacey, Andreas 204  
 Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and 205  
 Robin Rombach. SDXL: Improving Latent Diffusion Mod- 206  
 els for High-Resolution Image Synthesis. In *ICLR*, 2023. 1 207
- 208 [5] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya 208  
 Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, 209  
 Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen 210  
 Krueger, and Ilya Sutskever. Learning Transferable Visual 211

Table VI. Ablation study on the effect of inversion in SD1.5 across different datasets. Results are reported as mean and median scores; the best value for each metric is **bolded** and the second-best is underlined.

Dataset	SD1.5	Aesthetic [8]		PickScore [1]		HPS [9]		CLIP [5]	
		Mean	Median	Mean	Median	Mean	Median	Mean	Median
DrawBench	w inversion	<u>5.4230</u>	<u>5.4598</u>	<u>21.7540</u>	21.9768	<u>28.4336</u>	<u>28.6954</u>	<u>33.7922</u>	<u>34.5243</u>
	w/o inversion	<b>5.5822</b>	<b>5.6711</b>	<b>21.8486</b>	<b>21.9922</b>	<b>28.5449</b>	<b>28.7431</b>	<b>34.1385</b>	<b>34.8584</b>
PartiPrompts	w inversion	<u>5.5664</u>	<u>5.5970</u>	<u>21.8761</u>	<u>21.8947</u>	28.2564	<u>28.3230</u>	34.5021	<b>34.5227</b>
	w/o inversion	<b>5.7270</b>	<b>5.7612</b>	<b>22.0082</b>	<b>21.9777</b>	<b>28.3587</b>	<b>28.3920</b>	<b>34.5909</b>	<u>34.3746</u>
HPD v2	w inversion	<u>5.7916</u>	<u>5.7593</u>	<u>21.8716</u>	<u>21.8406</u>	<u>28.4611</u>	<u>28.4914</u>	<b>36.6221</b>	<b>36.9317</b>
	w/o inversion	<b>5.9261</b>	<b>5.9227</b>	<b>22.0227</b>	<b>21.9931</b>	<b>28.6379</b>	<b>28.6878</b>	<u>36.4029</u>	<u>36.7241</u>
Pick-a-Pic v2	w inversion	<u>5.6295</u>	<u>5.6464</u>	<u>21.4598</u>	<u>21.4744</u>	<u>27.6499</u>	<u>27.7286</u>	<b>34.7737</b>	<u>35.0402</u>
	w/o inversion	<b>5.7842</b>	<b>5.8168</b>	<b>21.5828</b>	<b>21.6021</b>	<b>27.8492</b>	<b>28.0742</b>	<u>34.6979</u>	<b>35.0903</b>

212 Models From Natural Language Supervision. In *ICML*,  
213 2021. 2, 3, 4, 5  
214 [6] Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Er-  
215 mon, Christopher D. Manning, and Chelsea Finn. Direct  
216 Preference Optimization: Your Language Model is Secretly  
217 a Reward Model. In *NeurIPS*, 2023. 1  
218 [7] Robin Rombach, Andreas Blattmann, Dominik Lorenz,  
219 Patrick Esser, and Björn Ommer. High-Resolution Image  
220 Synthesis with Latent Diffusion Models. In *CVPR*, 2022. 1  
221 [8] Christoph Schuhmann, Romain Beaumont, Richard Vencu,  
222 Cade Gordon, Ross Wightman, Mehdi Cherti, Theo  
223 Coombes, Aarush Katta, Clayton Mullis, et al. LAION-  
224 5B: An open large-scale dataset for training next generation  
225 image-text models. In *NeurIPS*, 2022. 2, 3, 4, 5  
226 [9] Xiaoshi Wu, Yiming Hao, Keqiang Sun, Yixiong Chen, Feng  
227 Zhu, Rui Zhao, and Hongsheng Li. Human Preference Score  
228 v2: A Solid Benchmark for Evaluating Human Preferences  
229 of Text-to-Image Synthesis. *arXiv.2306.09341*, 2023. 2, 3,  
230 4, 5  
231 [10] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gun-  
232 jan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yin-  
233 fei Yang, Burcu Karagol Ayan, et al. Scaling autoregressive  
234 models for content-rich text-to-image generation. In *TMLR*,  
235 2022. 2



User

You are an object detection system. Your task is as follows:

You will be provided with two images simultaneously. Your goals are:

Detect significant, identical object instances that are present in both images. (e.g., man, woman, dog, moon, landmark building)

You must follow these instructions:

Instance Matching: Carefully analyze both images to identify prominent and recognizable similar object instances to the greatest extent possible. The matching should follow these principles:

Priority for Identical Instances: First, seek to identify instances confirmed to be the same entity (e.g., the same person, the same specific pet, the same landmark building).

Relaxed Similarity Standards: If truly identical instances cannot be found, identify highly similar object instances, allowing for:

1. Minor variations in viewpoint, pose, or size.
2. Similar but not perfectly identical object categories (e.g., generic categories like "person", "car", "building").

Your output must be a clear list and strictly adhere to the following format:

Matched Instance: [Instance Category 1]

Matched Instance: [Instance Category 2]

Matched Instance: [Instance Category 3]



Planner

Matched Instance: [Eagle]



User

You are an object comparison system. Your task is as follows:

You will be provided with two cropped images simultaneously. Your goals are:

For each instance, evaluate and select the better-performing instance from the two images based on three dimensions.

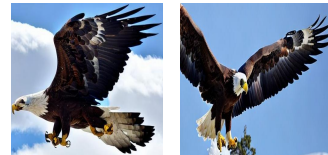
Quality Comparison: For each matched instance, evaluate the instance's performance in both images based on the following three dimensions and select the better one:

1. Anatomical/Structural Correctness: Are the object's proportions and structure correct and natural?
2. Naturalness & Artifacts: Does the object exhibit obvious AI-generated artifacts, blurriness, strange textures, or an unrealistic appearance?
3. Aesthetic Quality: Is the object visually appealing? Is the composition, color scheme, and overall concept harmonious and interesting?

Provide the final selection (Image A or Image B) without any detailed reason analysis.

Your output must be a clear list and strictly adhere to the following format:

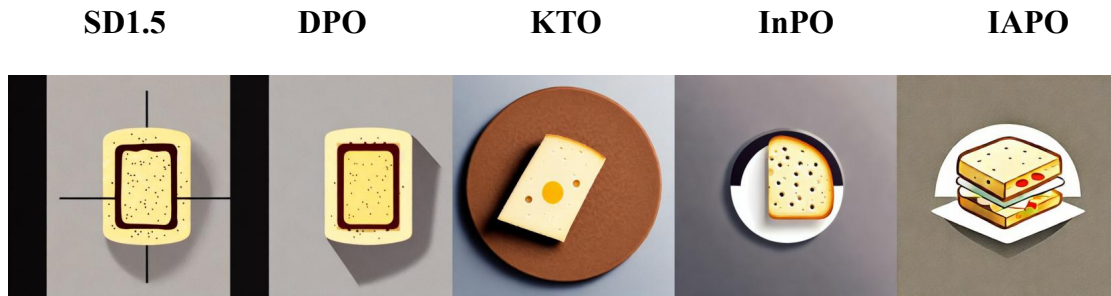
Quality Comparison Result: {Image A / Image B} is better



Judge

Quality Comparison Result: Image B is better

Figure II. Prompts of Planner and Judge.



A flat design illustration of a cheese sandwich with minimalistic line elements.



A girl with white hair and a school uniform, depicted in an illustration with warm clothes and a cold background.



A cartoon satanic priest depicted as an anthropomorphic lamb in a highly detailed 3D render.



Cartoon-style badger wearing a scarf against a green background.



The interior of a spaceship orbiting alpha centauri.

Figure III. Additional qualitative evaluation of IAPO-SD1.5 in comparison with Base-SD1.5, DPO-1.5, KTO-SD1.5 and InPO-SD1.5 on T2I generation tasks

SDXL

DPO

MaPO

InPO

IAPO



A white-haired girl in a pink sweater looks out a window in her bedroom.



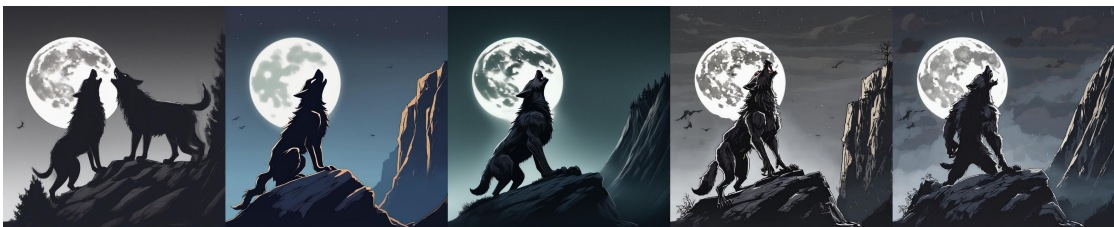
A lynx dressed in a flight suit.



A girl in a school uniform playing an electric guitar.



A sand monster amidst a tornado in the desert.



A werewolf howling on a cliff at night.

Figure IV. Additional qualitative evaluation of IAPO-SDXL in comparison with Base-SDXL, DPO-SDXL, MaPO-SDXL and InPO-SDXL on T2I generation tasks