

# DSH-BENCH: A DIFFICULTY- AND SCENARIO-AWARE BENCHMARK WITH HIERARCHICAL SUBJECT TAXONOMY FOR SUBJECT-DRIVEN TEXT-TO-IMAGE GENERATION

Anonymous authors

Paper under double-blind review

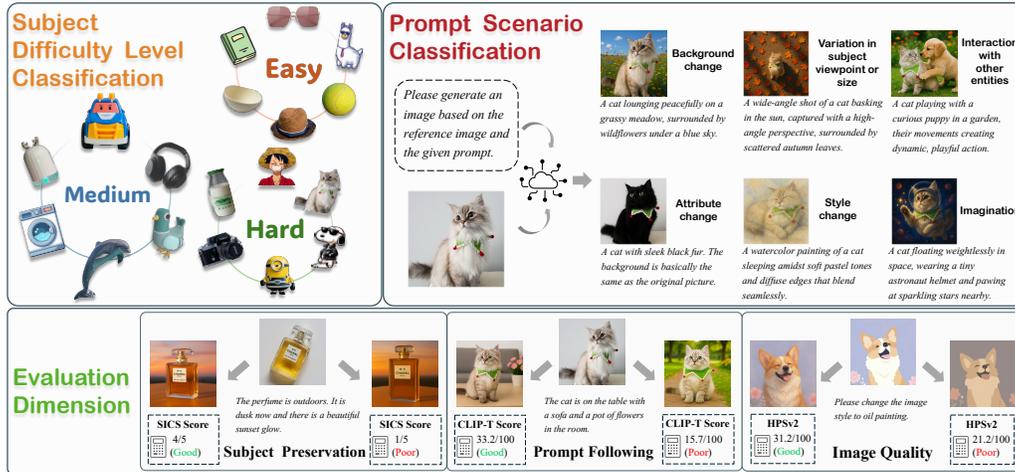


Figure 1: **Overview of DSH-Bench.** We curate a diverse dataset of subject images and categorize them into three difficulty levels—**easy**, **medium**, and **hard**—based on the complexity of preserving subject details. Leveraging GPT-4o’s capabilities, we systematically generate contextually appropriate prompts for various scenarios. The generated images are then rigorously evaluated across three key dimensions: **Subject Preservation**, **Prompt Following**, and **Image Quality**.

## ABSTRACT

Significant progress has been achieved in subject-driven text-to-image (T2I) generation, which aims to synthesize new images depicting target subjects according to user instructions. However, evaluating these models remains a significant challenge. Existing benchmarks exhibit critical limitations: 1) insufficient diversity and comprehensiveness in subject images, and 2) inadequate granularity in assessing model performance across different subject difficulty levels and prompt scenarios. To address these limitations, we propose DSH-Bench, a comprehensive benchmark that enables systematic multi-perspective analysis of subject-driven T2I models through three principal innovations: 1) a hierarchical taxonomy sampling mechanism ensuring comprehensive subject representation across 58 fine-grained categories, 2) an innovative classification scheme categorizing both subject difficulty level and prompt scenario for granular model capability assessment, and 3) a novel Subject Identity Consistency Score (SICS) metric demonstrating 9.4% higher correlation with human evaluation compared to existing measures in quantifying subject preservation. Through empirical evaluation of 15 subject-driven T2I models, DSH-Bench uncovers previously obscured limitations in current approaches while establishing concrete directions for future research.

---

# 1 INTRODUCTION

Subject-driven text-to-image (T2I) generation aims to generate images conditioned on both textual prompts and specific reference images. It has become feasible due to significant advancements in large-scale T2I generative models (Ding et al., 2021; Gafni et al., 2022; Saharia et al., 2022; Rombach et al., 2022a; Balaji et al., 2022; Chang et al., 2023; Kang et al., 2023; Dong et al., 2024). In subject-driven T2I generation, aside from image quality considerations, two other fundamental criteria must be satisfied: Subject Preservation and Prompt Following. Subject Preservation requires that the generated image maintain the details of the reference subject. Prompt Following demands that the generated image consistently reflects the content in the prompt. For example, a user might request an image of "his dog traveling around the world". In this scenario, the generated image must depict a dog identical to the reference image while illustrating the act of traveling as described.

Significant progress has been made in subject-driven T2I generation in recent years (Ruiz et al., 2023; Gal et al., 2022; Kumari et al., 2023; Wang et al., 2024a; Li et al., 2023a; Ye et al., 2023; Gal et al., 2023b; Wei et al., 2023; Hu et al., 2024b; Qiu et al., 2023). One approach involves fine-tuning general T2I models to create specialized models that reproduce specific subjects present in the training datasets. Alternatively, encoder-based methods achieve subject preservation by adapting features to incorporate reference subject into a general T2I model. Despite these advancements, challenges remain in comprehensively and effectively evaluating the actual performance of these models. An effective evaluation method should not only provide a comprehensive and unbiased assessment, but also align with human perception to ensure reliable measurement. Furthermore, the evaluation method is expected to provide valuable insights for future research. However, current benchmarks (Ruiz et al., 2023; Kumari et al., 2023; Chen et al., 2023a; Wang et al., 2024b; Peng et al., 2025) are limited by insufficient diversity and comprehensiveness in subject image collection, which restricts the thoroughness of model evaluation. In addition, they do not facilitate a detailed understanding of subject difficulty and prompt scenarios, thus constraining the depth of insights obtainable from the evaluation. As shown in Figure 3, our analysis of numerous model-generated instances reveals that different subject images and prompts place varying demands on a model’s ability. For example, although subject-driven T2I models are capable of effectively preserving the details of relatively simple objects (e.g., a tennis ball), they often struggle to accurately reproduce objects with more intricate features (e.g., a camera). This observation highlights the importance of categorizing the subject difficulty and prompt scenario to better assess model performance. To address these requirements, we introduce DSH-Bench, a novel benchmark offers three notable advantages:

1. **The diversity of subject images in DSH-Bench is substantially greater** To mitigate evaluation bias caused by low diversity of subject images, we employ a hierarchical taxonomy in image collection. We referenced COCO (Lin et al., 2014) and ImageNet (Deng et al., 2009) in the hierarchical taxonomy construction. As shown in Figure 4(a), the widely used DreamBench includes only 6 categories and 30 subjects. In contrast, our benchmark expands the dataset to 48 categories and 459 subjects—representing an increase of  $8\times$  and  $15\times$ , respectively. Although DreamBench++ (Peng et al., 2025) offers 150 subjects, its diversity is constrained by its image collection. Notably, 33% of our categories are not represented in DreamBench++. Therefore, benefiting from DSH-Bench’s greater subject diversity, we enable more comprehensive evaluation of models.

2. **An innovative classification scheme for subject difficulty level and prompt scenario** Figure 3 shows the model’s performance varies significantly with different samples, highlighting the necessity for a classification of both subject image and prompt. Although DreamBench++ categorizes prompts based on their perceived difficulty, the criteria underlying this classification are not clearly defined. Additionally, DreamBench++ does not analyze the difficulty levels associated with different subjects. To address these limitations, we categorize subjects into three difficulty levels (easy, medium, and hard) according to the difficulty of preserving visual appearance and classify prompts into six scenarios (background change, variation in subject viewpoint or size, interaction with other entities, attribute change, style change, imagination). As a result, our approach enables a more comprehensive and granular analysis of the challenges faced by current models.

3. **A human-aligned and more efficient metric for subject preservation** DreamBench++ replaces CLIP (Radford et al., 2021) and DINO (Caron et al., 2021) with GPT-4o (OpenAI, 2024) for evaluation, resulting in improved alignment with human evaluation. However, our benchmark reveals that per-model evaluation under this paradigm requires approximately 20,000 API calls to

108  
109  
110  
111  
112  
113  
114  
115  
116  
117  
118  
119  
120  
121  
122  
123  
124  
125  
126  
127  
128  
129  
130  
131  
132  
133  
134  
135  
136  
137  
138  
139  
140  
141  
142  
143  
144  
145  
146  
147  
148  
149  
150  
151  
152  
153  
154  
155  
156  
157  
158  
159  
160  
161



Figure 2: Qualitative comparison of subject preservation between SICS and the other methods.

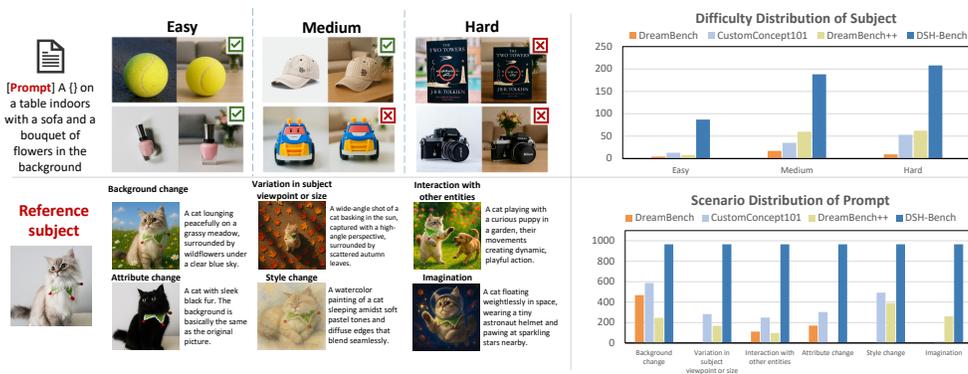


Figure 3: Qualitative comparison of generated images under different difficulty levels and scenarios.

GPT-4o, incurring prohibitive computational costs exceeding \$400 for each evaluation. To address the limitation, we introduce **Subject Identity Consistency Score (SICS)**, which innovatively focuses on subject-level consistency rather than merely relying on embedding comparisons. Firstly, five annotators label a training dataset containing 5,000 image-text pairs, focusing on subject preservation evaluation. We then fine-tune Qwen2.5-VL-7B (Bai et al., 2025) on this dataset, which leads the model to focus on core visual attributes rather than high-level semantics. Finally, we use Kendall’s  $\tau$  value to quantify the alignment between model outputs and human evaluation. Experimental results demonstrate that SICS achieves a statistically significant improvement, outperforming DreamBench++ by 9.4% in human evaluation correlation metrics. Figure 2 presents a partial qualitative comparison of concept preservation between SICS and the other assessment methods.

**Takeaways** We present some insightful findings from evaluating fifteen methods: i) Our evaluation reveals that no single method demonstrates consistently robust performance across all categories. Therefore, implementing hierarchical taxonomy sampling of subject images is critical for mitigating potential evaluation biases. ii) All methods exhibit degraded performance on hard subject images. It is crucial to enhance models’ ability to encode and reconstruct complex subject details more effectively in future research. iii) The subject-driven T2I model’s capability for different prompt scenarios is not robust. Future research on subject-driven T2I generation should focus on optimizing for adaptation to a variety of prompt scenarios.

In summary, our contributions are as follows: 1) We employ a hierarchical taxonomy in image collection to ensure both the diversity and comprehensiveness of subject images. 2) We propose an innovative classification scheme to categorize subject difficulty levels and prompt scenarios. This scheme enables us to obtain valuable insights. 3) We propose a human-aligned metric to evaluate subject preservation, which offers greater efficiency compared to DreamBench++. We are open-sourcing DSH-Bench, including subject images, prompts, generated images and related code.

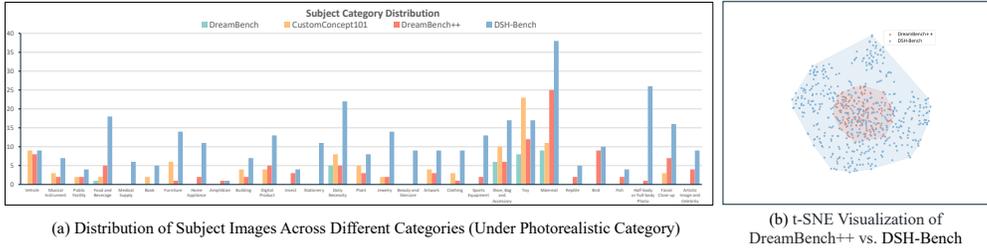


Figure 4: **Distribution of subject images.** (a) Category-wise image distribution for our benchmark versus prior benchmarks. (b) t-SNE comparison of images between DSH-Bench and DreamBench++.

## 2 DSH-BENCH

This section provides an overview of the primary components of DSH-Bench. Section 2.1 outlines the data construction process. Section 2.2 introduces the definitions and evaluation methods for three evaluation dimensions. *A detailed explanation is available in the supplementary materials.*

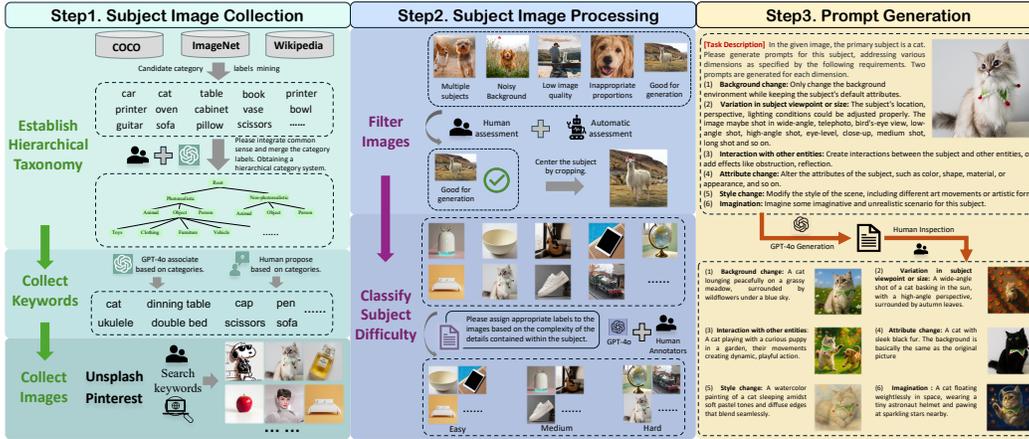


Figure 5: **Dataset construction process of DSH-Bench.** We construct a hierarchical taxonomy to obtain a comprehensive set of keywords. Then we collect web images using these keywords. After performing both manual review and automated filtering of the images, we classify the difficulty of subject images and use GPT-4o to generate prompts for each subject image.

### 2.1 BENCHMARK DATASET CONSTRUCTION

#### 2.1.1 SUBJECT IMAGE COLLECTION

**Hierarchical Taxonomy Establishment** As shown in Figure 5, we establish a hierarchical taxonomy. For the first- and second-level categories, we primarily refer to existing benchmarks from prior studies (Ruiz et al., 2023; Kumari et al., 2023; Peng et al., 2025), resulting in two first-level categories and six second-level categories. For the third-level categories, we first reference COCO and ImageNet to compile a list of candidate category labels, then utilize GPT-4o to consolidate them into 58 refined categories. The final hierarchical taxonomy is confirmed and refined through co-authors’ discussion. The detailed process and the category contents are provided in Appendix A.

**Keyword Collection & Internet Image Collection** In DreamBench++, keywords collection relies on GPT-4o and human input. The approach does not adequately ensure the diversity of the obtained keywords, potentially introducing bias during the image collection process. In contrast, DSH-Bench derives keywords from a hierarchical taxonomy. For each third-level category, we use GPT-4o to generate associated keywords, which are further supplemented by humans. All keywords are

then consolidated and deduplicated, resulting in a final set of **400** unique keywords—surpassing DreamBench++’s 300. The specific keywords are provided in the Appendix B. Given a set of selected keywords, we retrieve images from Unsplash ([uns](#)) and Pinterest ([pin](#)). Keywords without suitable images are discarded. *Each image’s copyright status has been verified for academic suitability.*

### 2.1.2 SUBJECT IMAGE PROCESSING

**Image Filtering** To filter unsuitable images, human annotators remove images with multiple subjects and noisy backgrounds. We use aesthetic score (Xu et al., 2024) and SAM (Kirillov et al., 2023) to filter images with low image quality and inappropriate proportions of subject regions. The curated images are subsequently cropped to centralize the reference subject.

**Subject Difficulty Level Classification** As illustrated in Figure 3, the model’s performance varies considerably across different samples. To derive meaningful insights, we classify the subject images according to the difficulty level that the model experiences in preserving details of the reference subject. We define three subject difficulty levels, including (1) **Easy**: Subjects characterized by minimal surface complexity and homogeneous textural properties, exemplified by smooth-surfaced objects such as a ceramic mug with uniform coloration. These cases present negligible challenges for detail preservation due to their structural regularity. (2) **Medium**: Subjects containing discernible high-frequency features while maintaining global structural coherence, such as cylindrical containers with legible typographic elements. These cases require intermediate detail preservation capabilities. (3) **Hard**: Subjects exhibiting non-uniform texture distributions and multi-scale geometric details, typified by objects like book covers containing fine-grained calligraphic elements. Such cases expose model limitations in maintaining structural fidelity and textural granularity under complex topological constraints. We utilize GPT-4o to classify the subject images according to the aforementioned criteria. Subsequently, all images are reviewed by five human annotators to ensure accuracy and consistency.

### 2.1.3 PROMPT GENERATION

Although DreamBench++ categorizes prompts based on their perceived difficulty, it does not provide empirical evidence to substantiate the criterion. To address this limitation, we organize the prompts according to specific application scenarios, dividing them into six categories, including (1) **Background change (BC)**: scenarios involving changes in background elements. (2) **Variation in subject viewpoint or size (VS)**: scenarios that entail changes in camera angle, which may include variations in subject size, lighting, or shadows. (3) **Interaction with other entities (IE)**: scenarios requiring complex interactions with additional entities, potentially resulting in occlusion and necessitating adherence to physical plausibility. (4) **Attribute change (AC)**: scenarios involving modifications to certain attributes of the subject, such as color or shape. (5) **Style change (SC)**: scenarios involving alterations in the artistic or visual style of the subject. (6) **Imagination (IM)**: scenarios where the target image depicts an imagined or fictional scene. We generate two prompts for each scenario. The specific instructions are depicted in Figure 5. All prompts are reviewed by five human annotators to ensure they are ethical and free from defects. For the specific verification procedure, please refer to the Appendix E.3. Finally, we obtain a total of **459** high-quality images and **5,508** prompts. Figure 3 shows the distribution of subject image difficulty levels and prompt scenarios. We visualize the t-SNE of images from our benchmark and DreamBench++ in Figure 4(b). The results indicate that our benchmark achieves superior diversity.

## 2.2 EVALUATION DIMENSION

Previous notable works (Ruiz et al., 2023; Gal et al., 2022; Kumari et al., 2023; Wang et al., 2024a) evaluate the performance of subject-driven T2I models from two perspectives: Subject Preservation and Prompt Following. Mao et al. (2024) also uses ImageReward (Xu et al., 2023) to evaluate image quality. Therefore, DSH-Bench evaluates from the three aforementioned dimensions.

**Subject Preservation** DreamBench++ utilizes GPT-4o for evaluation to improve alignment with human assessments. However, the GPT-4o-based method is prohibitively expensive. To address this limitation, we propose a novel metric—**Subject Identity Consistency Score (SICS)**. Firstly, we establish a scoring criterion for assessing subject preservation, the details are provided in Appendix E.2. Five annotators label the collected image pairs according to the criterion. During the annotation process, each image pair is not only assigned a score but also accompanied by an explanation. Previous

work (Wei et al., 2022) has indicated that labeled data with explanatory reasoning can help models better understand the underlying logic and reasoning behind the labels. We then perform meticulous fine-tuning of the model using this annotated dataset. During fine-tuning, SICS leverages prompts to explicitly prioritize subject consistency rather than global semantics, mitigating background and style artifacts that commonly bias CLIP-based approaches and yielding closer alignment with the goals of subject-consistency evaluation. Although GPT-4o demonstrates outstanding performance across a wide range of tasks, it has not been specifically optimized for subject preservation evaluation. More details of the SICS metric can be found in Appendix E.2.

**Prompt Following** Prompt following primarily evaluates whether a model can generate images that accurately correspond to textual prompts. DreamBench++ has demonstrated that the CLIP-T score is highly consistent with human annotations. Therefore, we also adopt CLIP-T score as the evaluation metric for prompt following.

**Image Quality** HPSv2 (Wu et al., 2023) utilizes professionally annotated data to more accurately reflect human aesthetic preferences for generated images. Previous studies (Sun et al., 2025) demonstrate that models optimized with HPSv2 achieve superior performance in image quality assessment compared to existing approaches. Therefore, we adopt HPSv2 for image quality evaluation.

### 3 EXPERIMENT

#### 3.1 EXPERIMENT SETUP

**Implementation Details** We conduct experiments on two mainstream approaches: *i) Finetuning-based:* 1) Textual Inversion(TI) (Gal et al., 2023a), 2) DreamBooth, 3) Custom Diffusion, 4) Hiper (Han et al., 2023), 5) NeTI (Alaluf et al., 2023). *ii) Encoder-based:* 1) BLIP-Diffusion (Li et al., 2023a), 2) IP-Adapter (Ye et al., 2023), 3) MS-Diffusion (Wang et al., 2024b), 4) Emu2 (Sun et al., 2024), 5) OminiControl (Tan et al., 2024), 6) SSR-Encoder (Zhang et al., 2024), 7) RealCustom++ (Mao et al., 2024), 8) OmniGen (Xiao et al., 2024), 9)  $\lambda$ -Eclipse (Patel et al., 2024), 10) UNO (Wu et al., 2025). Our experiments are conducted using the official implementations to guarantee reliability and fairness. More details can be found in Appendix E.

**Human Annotation** All annotation tasks, including labeling of the SICS training datasets, were conducted by the same five human annotators. We provide the annotators with detailed labeling guidelines and sufficient training to ensure they fully understand the subject-driven T2I generation task and could provide unbiased and discriminative scores. For additional details regarding the human annotation process, please see the Appendix E.4.

Table 1: The human alignment degree among different metrics, measured by **Kendall’s  $\tau$  value** and **Spearman correlation coefficient value**. H: Human, G: GPT-4o, D: DINO, Dv2: DINOv2, CB: CLIP-B, CL: CLIP-L, S: SICS. Bold font is used to denote the maximum value in a row.

Method	Kendall $\uparrow$					Spearman $\uparrow$						
	H-CB	H-CL	H-D	H-Dv2	H-G	H-S	H-CB	H-CL	H-D	H-Dv2	H-G	H-S
BLIP-Diffusion	0.228	0.176	0.285	0.167	<u>0.354</u>	<b>0.531</b>	0.285	0.215	0.350	0.206	0.383	<b>0.554</b>
IP-Adapter	0.294	0.296	0.258	0.290	<u>0.419</u>	<b>0.622</b>	0.364	0.371	0.325	0.364	<u>0.459</u>	<b>0.657</b>
MS-Diffusion	0.158	0.090	0.116	0.122	0.119	<b>0.178</b>	<b>0.194</b>	0.109	0.144	0.156	0.131	0.189
OminiControl	<u>0.375</u>	0.371	0.337	0.348	0.650	<b>0.713</b>	0.490	0.486	0.441	0.453	0.729	<b>0.764</b>
SSR-Encoder	0.264	0.338	0.295	0.348	<u>0.504</u>	<b>0.664</b>	0.328	0.421	0.368	0.434	<u>0.549</u>	<b>0.697</b>
UNO	0.249	0.218	<u>0.299</u>	0.240	<u>0.236</u>	<b>0.385</b>	0.340	0.297	<u>0.390</u>	0.312	0.268	<b>0.426</b>
RealCustom++	0.181	0.128	<u>0.206</u>	0.241	<u>0.291</u>	<b>0.464</b>	0.229	0.162	0.266	0.303	<u>0.325</u>	<b>0.511</b>
OmniGen	0.465	0.396	0.344	0.349	<u>0.617</u>	<b>0.621</b>	0.579	0.497	0.440	0.456	<b>0.697</b>	0.667
$\lambda$ -Eclipse	0.143	0.233	0.084	0.103	<u>0.325</u>	<b>0.375</b>	0.176	0.287	0.103	0.127	<u>0.352</u>	<b>0.393</b>
Custom Diffusion	0.316	0.336	0.382	0.425	<u>0.487</u>	<b>0.642</b>	0.388	0.409	0.470	0.519	<u>0.512</u>	<b>0.654</b>
DreamBooth	0.639	0.591	0.537	0.429	<u>0.647</u>	<b>0.692</b>	0.733	0.721	0.661	0.537	0.705	<b>0.740</b>
Textual Inversion	0.482	0.459	0.447	0.438	0.541	<b>0.568</b>	<u>0.587</u>	0.559	0.545	0.534	0.582	<b>0.590</b>
HiPer	0.338	0.387	0.351	0.404	<u>0.584</u>	<b>0.625</b>	0.417	0.469	0.430	0.496	0.629	<b>0.655</b>
NeTI	0.469	0.456	0.431	0.417	<u>0.617</u>	<b>0.728</b>	0.573	0.561	0.529	0.512	<u>0.682</u>	<b>0.778</b>
ALL	0.416	0.411	0.350	0.376	<u>0.619</u>	<b>0.677</b>	0.529	0.522	0.451	0.483	<u>0.697</u>	<b>0.734</b>

#### 3.2 MAIN RESULTS

**SICS Results** Table 1 presents a rigorous study of human alignment using *Kendall’s  $\tau$  value* (KDV) and *Spearman correlation coefficient value* (SCV) (metric selection rationale in Appendix E.2). Our experimental results demonstrate that **SICS achieves superior alignment with human evaluations**

Table 2: **Evaluation of Subject-driven T2I generation.** DB: DreamBench, DB++: DreamBench++, HB: DSH-Bench. All scores are normalized to 0-1. Bold indicates the minimum value in each row for a given evaluation dimension..

Method	Subject Preservation			Prompt Following			Image Quality		
	DB	DB++	HB	DB	DB++	HB	DB	DB++	HB
BLIP-Diffusion	0.229	0.216	<b>0.204</b>	0.291	0.278	<b>0.277</b>	0.267	0.254	<b>0.223</b>
IP-Adapter	0.230	0.244	<b>0.229</b>	0.321	0.318	<b>0.315</b>	0.291	0.296	<b>0.266</b>
MS-Diffusion	<b>0.316</b>	0.346	0.352	<b>0.332</b>	0.339	0.338	0.311	0.314	<b>0.294</b>
OminiControl	0.279	0.268	<b>0.258</b>	<b>0.325</b>	0.337	0.334	0.312	0.308	<b>0.290</b>
SSR-Encoder	0.231	<b>0.202</b>	<b>0.202</b>	0.290	<b>0.287</b>	0.295	0.273	0.270	<b>0.247</b>
UNO	<b>0.409</b>	0.410	<b>0.409</b>	<b>0.317</b>	0.322	0.323	0.304	0.297	<b>0.278</b>
Emu2	0.360	0.343	<b>0.341</b>	<b>0.291</b>	0.309	0.304	0.272	0.278	<b>0.260</b>
RealCustom++	0.377	0.380	<b>0.375</b>	<b>0.325</b>	0.329	0.332	0.316	0.314	<b>0.298</b>

Table 3: **DSH-Bench leaderboard.** The models are ranked by the final score  $S_h$ . We only present the top models; the complete ranking can be found in the Appendix D.2.

Method	T2I Model	Subject Preservation	Prompt Following	Image Quality	$S_h \uparrow$
UNO	FLUX.1-dev	<b>0.409</b>	0.323	0.278	<b>0.252</b>
RealCustom++	SDXL	0.375	0.332	<b>0.294</b>	0.251
MS-Diffusion	SDXL	0.352	<b>0.338</b>	<b>0.294</b>	0.248
Emu2	SDXL	0.341	0.304	0.260	0.228
OminiControl	FLUX.1-schnell	0.258	0.334	0.290	0.218
IP-Adapter	SDXL	0.256	0.292	0.266	0.199
$\lambda$ -Eclipse	SDXL	0.229	0.315	0.242	0.198
OmniGen	SD v1.5	0.202	0.295	0.265	0.183
SSR-Encoder	SDXL	0.188	0.322	0.247	0.181
NeTI	SD v1.4	0.192	0.301	0.234	0.176
BLIP-Diffusion	SD v1.5	0.204	0.277	0.223	0.174
DreamBooth	SD v1.5	0.158	0.321	0.245	0.164
HiPer	SD v1.4	0.135	0.318	0.247	0.151
Textual Inversion	SD v1.5	0.109	0.299	0.225	0.129
Custom Diffusion	SD v1.4	0.062	0.323	0.240	0.091

compared to existing methods, showing consistently higher agreement across both correlation metrics in most experimental settings. Although SICS attains second-highest correlation scores in MS-Diffusion and OmniGen, it significantly outperforms GPT-4o (*GPT-4o refers to the evaluation method used in DreamBench++*) by **9.37%** (KDV) and **5.31%** (SCV). This performance gap strongly suggests SICS’s enhanced capability in modeling human evaluation. Notably, GPT-4o exhibits greater consistency with human evaluation than CLIP and DINO, aligning with DreamBench++ findings. Importantly, our proposed SICS metric surpasses all existing metrics in human judgment consistency.

**Quantitative & Qualitative Results** Table 2 shows overall evaluation results. The results show that: **i) DSH-Bench poses more significant challenges than existing benchmarks.** For subject preservation and image quality, the majority of methods consistently yield lower scores on DSH-Bench. The result can be attributed to the hierarchical taxonomy sampling method employed, which allows our dataset to more accurately represent the true data distribution. Moreover, it highlights that benchmarks derived from true distributions present greater challenges. **ii)** For prompt following, DreamBench yields slightly lower scores than DSH-Bench for certain methods. In DreamBench, prompts requiring attribute change constitute 22.7%, which is higher than the 16.7% observed in DSH-Bench. Figure 7(b) indicates that all methods exhibit relatively poor average performance on prompts involving attribute change. **iii)** Table 3 shows that there exists a trade-off between subject preservation and prompt following. We plot the Pareto frontier (see in Appendix D.1) using the data presented in Table 3. The primary objective is to identify a Pareto optimal solution that effectively balances the two objectives. *Additional results and discussions can be found in Appendix D.2.*

**Leaderboard** In order to assess a model’s overall capability, we define the final score as:

$$S_h = \frac{3}{\frac{\lambda}{SP} + \frac{\gamma}{PF} + \frac{\mu}{IQ}} \quad (1)$$

SP, PF, and IQ represent the scores for Subject Preservation, Prompt Following, and Image Quality, respectively.  $\lambda, \gamma, \mu$  are the weights assigned to the importance of each corresponding dimension. In this study, we set  $\lambda = 1.5, \gamma = 1.5, \mu = 1$ , as subject preservation and prompt following are of paramount importance in subject-driven T2I generation. The harmonic mean requires strong

performance across all dimensions to yield a high overall score. We rank models by  $S_h$  scores in Table 3. UNO exhibits relatively strong overall performance. We attribute this improvement to the novel architectural design of UNO and the minimal yet effective modifications implemented in DiT.

## 4 ANALYSIS

In this section, we conduct a detailed analysis of the performance of all methods based on the hierarchical category, the subject difficulty level classification, and the prompt scenario classification:

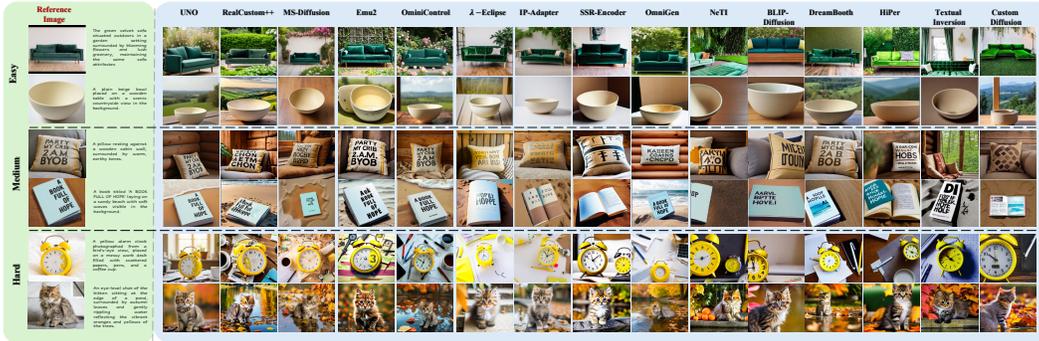


Figure 6: Examples generated by methods listed in the leaderboard. Best viewed when zoomed in.

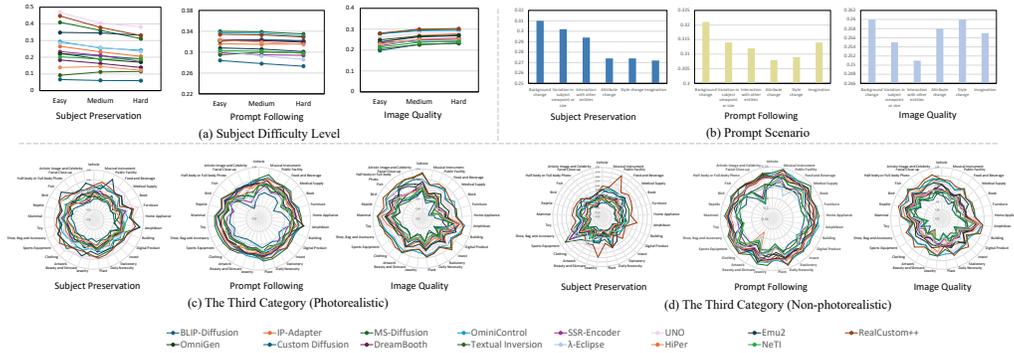


Figure 7: Comparison for DSH-Bench scores in different evaluation dimensions. The specific metric values are provided in the Appendix D.2. Best viewed when zoomed in.

**A scientific and comprehensive subject image sampling method is necessary** Figure 7(c) and Figure 7(d) present the performance of various methods in the third-level categories. The results reveal that model robustness varies considerably among categories. For example, performance in categories "artwork" (both photorealistic and non-photorealistic) is substantially lower. This disparity suggests that the absence of subject images from specific categories can lead to biased evaluation results, highlighting the importance of data diversity. Furthermore, Figure 7 also demonstrates that none of the current models perform well across all categories. We hypothesize that this may be related to the varying complexity of the subjects within different categories. A more detailed analysis of model performance in different categories can be found in Appendix D.1.

**Current subject-driven T2I models exhibit performance degradation on hard level subjects** As illustrated in Figure 7(a), the model exhibits substantial variation in performance across different difficulty levels: 1) For subject preservation, there is a pronounced decline in performance as the difficulty of the subject images increases. The model achieves significantly better results on images classified as simple compared to those categorized as hard. This observation supports the validity of our image difficulty classification scheme. 2) For prompt following, Figure 7(a) shows that the capability of the models is minimally influenced by the subject difficulty level. This could be explained by the fact that CLIP-T primarily emphasizes overall semantic information. Consequently,

---

432 as long as the generated image correctly represents the general category and overall shape, the  
433 evaluation score is unlikely to be substantially reduced, even if finer details are not perfectly captured.  
434 *Given these findings, it is crucial to enhance models' ability to encode and reconstruct complex*  
435 *subject details more effectively in future research endeavors.*

436 **The subject-driven T2I capability for different prompt scenarios is not robust** Figure 7(b)  
437 shows the average performance of all models across six prompt scenarios. The results show that: 1)  
438 In BC, VS, and IE scenarios, the model's performance consistently declines across all evaluation  
439 dimensions. This trend suggests that the difficulty of the scenarios increases progressively from  
440 BC to IE. Notably, the finding that the IE scenario is more challenging than the BC scenario aligns  
441 with intuitive expectations. 2) For subject preservation, the model's average performance across  
442 the AC, SC, and IM prompt scenarios remains relatively low. This could be because the generated  
443 subjects undergo partial modifications relative to the original subjects in these three scenarios. *Given*  
444 *these findings, more emphasis should be placed on enhancing methods for IE prompt scenario. For*  
445 *instance, increasing the volume of training data tailored to these specific contexts.*

## 447 5 RELATED WORK

### 449 5.1 SUBJECT-DRIVEN TEXT-TO-IMAGE GENERATION

451 In recent years, subject-driven T2I generation has attracted significant research attention (Ruiz  
452 et al., 2023; Gal et al., 2022; Kumari et al., 2023; Wang et al., 2024a; Li et al., 2023a; Gal et al.,  
453 2023b;a; Wei et al., 2023; Hu et al., 2024b; Qiu et al., 2023). Within the context of diffusion models,  
454 optimization-based model (Voynov et al., 2023; Liu et al., 2023; Hua et al., 2023; Hao et al., 2023)  
455 enables subject-driven generation by introducing lightweight parameters and performs parameter-  
456 efficient fine-tuning for each subject. In contrast, the encoder-based methods (Shi et al., 2023; Ma  
457 et al., 2024; Chen et al., 2023b; Li et al., 2023b; Le et al., 2024; Rowles et al., 2024; Zeng et al., 2024;  
458 Hu et al., 2024a; Huang et al., 2025a; Xiong et al., 2025; Patashnik et al., 2025; Wu et al., 2025;  
459 Huang et al., 2025b; He et al., 2025) leverage additional image encoders and network layers to encode  
460 the reference image of the subject. IP-Adapter (Ye et al., 2023) introduces cross-attention through an  
461 additional image encoder to incorporate control signals. Furthermore, SSR-Encoder (Zhang et al.,  
462 2024) enhances identity preservation without necessitating further fine-tuning when introducing new  
463 concepts. The Diffusion Transformers (Peebles & Xie, 2023; Podell et al., 2023; Rombach et al.,  
464 2022b) uses transformer as a denoising network to iteratively refine noisy image tokens. Based on  
465 these foundation models, approaches like OminiControl (Tan et al., 2024) and UNO (Wu et al., 2025)  
466 explore the inherent image reference capabilities of transformers.

### 467 5.2 SUBJECT-DRIVEN T2I GENERATION BENCHMARK

469 Evaluation of subject-driven T2I relies on diverse metrics. For image quality, several notable  
470 studies (Xu et al., 2023; Kirstain et al., 2023; Wu et al., 2023; Alaluf et al., 2023; Xu et al., 2024;  
471 Wang et al., 2025) have been proposed. Subject preservation is typically measured by learning-based  
472 metrics that compare deep feature distances, often using embeddings from large vision models such as  
473 CLIP (Radford et al., 2021) and DINO (Caron et al., 2021), as well as image-retrieval scores (Liu et al.,  
474 2021). To better align with human perception, DreamSim (Fu et al., 2023) emphasizes foreground  
475 objects when assessing image similarity. Semantic consistency is commonly measured with the CLIP  
476 score. DreamBench lacks diversity in subjects and prompts, and although DreamBench++ expands to  
477 150 subjects, it still lacks systematic categorization, hindering meaningful analysis.

## 479 6 CONCLUSION

481 This paper introduces a novel benchmark called DSH-Bench, designed specifically for subject-driven  
482 T2I generation. Key features include: 1) a hierarchical category system in image collection to ensure  
483 both the diversity and comprehensiveness of subject images; 2) an innovative classification scheme  
484 for categorizing subject difficulty levels and prompt scenarios to obtain valuable insights; and 3) a  
485 human-aligned and more efficient metric for subject preservation. The benchmark will be publicly  
available to support the advancement in the subject-driven T2I generation era.

486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539

---

## REFERENCES

- Yuval Alaluf, Elad Richardson, Gal Metzer, and Daniel Cohen-Or. A neural space-time representation for text-to-image personalization. *ACM Transactions on Graphics (TOG)*, 42(6):1–10, 2023.
- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng, Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. Qwen2.5-vl technical report, 2025. URL <https://arxiv.org/abs/2502.13923>.
- Yogesh Balaji, Seungjun Nah, Xun Huang, Arash Vahdat, Jiaming Song, Qinsheng Zhang, Karsten Kreis, Miika Aittala, Timo Aila, Samuli Laine, et al. ediff-i: Text-to-image diffusion models with an ensemble of expert denoisers. *arXiv preprint arXiv:2211.01324*, 2022.
- Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 9650–9660, 2021.
- Huiwen Chang, Han Zhang, Jarred Barber, Aaron Maschinot, Jose Lezama, Lu Jiang, Ming-Hsuan Yang, Kevin Patrick Murphy, William T Freeman, Michael Rubinstein, et al. Muse: Text-to-image generation via masked generative transformers. In *International Conference on Machine Learning*, pp. 4055–4075. PMLR, 2023.
- Wenhu Chen, Hexiang Hu, Yandong Li, Nataniel Ruiz, Xuhui Jia, Ming-Wei Chang, and William W Cohen. Subject-driven text-to-image generation via apprenticeship learning. *Advances in Neural Information Processing Systems*, 36:30286–30305, 2023a.
- Zhuowei Chen, Shancheng Fang, Wei Liu, Qian He, Mengqi Huang, Yongdong Zhang, and Zhendong Mao. Dreamidentity: Improved editability for efficient face-identity preserved image generation, 2023b. URL <https://arxiv.org/abs/2307.00300>.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.
- Ming Ding, Zhuoyi Yang, Wenyi Hong, Wendi Zheng, Chang Zhou, Da Yin, Junyang Lin, Xu Zou, Zhou Shao, Hongxia Yang, et al. Cogview: Mastering text-to-image generation via transformers. *Advances in neural information processing systems*, 34:19822–19835, 2021.
- Runpei Dong, Chunrui Han, Yuang Peng, Zekun Qi, Zheng Ge, Jinrong Yang, Liang Zhao, Jianjian Sun, Hongyu Zhou, Haoran Wei, et al. Dreamllm: Synergistic multimodal comprehension and creation. In *ICLR*, 2024.
- Stephanie Fu, Netanel Tamir, Shobhita Sundaram, Lucy Chai, Richard Zhang, Tali Dekel, and Phillip Isola. Dreamsim: Learning new dimensions of human visual similarity using synthetic data. In *Advances in Neural Information Processing Systems*, volume 36, pp. 50742–50768, 2023.
- Oran Gafni, Adam Polyak, Oron Ashual, Shelly Sheynin, Devi Parikh, and Yaniv Taigman. Make-a-scene: Scene-based text-to-image generation with human priors. In *European Conference on Computer Vision*, pp. 89–106. Springer, 2022.
- Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H. Bermano, Gal Chechik, and Daniel Cohen-Or. An image is worth one word: Personalizing text-to-image generation using textual inversion, 2022. URL <https://arxiv.org/abs/2208.01618>.
- Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit Haim Bermano, Gal Chechik, and Daniel Cohen-or. An image is worth one word: Personalizing text-to-image generation using textual inversion. In *The Eleventh International Conference on Learning Representations*, 2023a. URL <https://openreview.net/forum?id=NAQvF08TcyG>.
- Rinon Gal, Moab Arar, Yuval Atzmon, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or. Encoder-based domain tuning for fast personalization of text-to-image models. *ACM Transactions on Graphics (TOG)*, 42(4):1–13, 2023b.

---

540 Inhwa Han, Serin Yang, Taesung Kwon, and Jong Chul Ye. Highly personalized text embedding for  
541 image manipulation by stable diffusion. *arXiv preprint arXiv:2303.08767*, 2023.

542

543 Shaozhe Hao, Kai Han, Shihao Zhao, and Kwan-Yee K Wong. Vico: Plug-and-play visual condition  
544 for personalized text-to-image generation. *arXiv preprint arXiv:2306.00971*, 2023.

545

546 Junjie He, Yuxiang Tuo, Binghui Chen, Chongyang Zhong, Yifeng Geng, and Liefeng Bo. Anystory:  
547 Towards unified single and multiple subject personalization in text-to-image generation, 2025.  
548 URL <https://arxiv.org/abs/2501.09503>.

549 Hexiang Hu, Kelvin C. K. Chan, Yu-Chuan Su, Wenhui Chen, Yandong Li, Kihyuk Sohn, Yang  
550 Zhao, Xue Ben, Boqing Gong, William Cohen, Ming-Wei Chang, and Xuhui Jia. Instruct-imagen:  
551 Image generation with multi-modal instruction, 2024a. URL [https://arxiv.org/abs/  
552 2401.01952](https://arxiv.org/abs/2401.01952).

553 Hexiang Hu, Kelvin CK Chan, Yu-Chuan Su, Wenhui Chen, Yandong Li, Kihyuk Sohn, Yang Zhao,  
554 Xue Ben, Boqing Gong, William Cohen, et al. Instruct-imagen: Image generation with multi-  
555 modal instruction. In *Proceedings of the IEEE/CVF conference on computer vision and pattern  
556 recognition*, pp. 4754–4763, 2024b.

557

558 Miao Hua, Jiawei Liu, Fei Ding, Wei Liu, Jie Wu, and Qian He. Dreamtuner: Single image is enough  
559 for subject-driven generation, 2023. URL <https://arxiv.org/abs/2312.13691>.

560

561 Linyan Huang, Haonan Lin, Yanning Zhou, and Kaiwen Xiao. Flexip: Dynamic control of preser-  
562 vation and personality for customized image generation, 2025a. URL [https://arxiv.org/  
563 abs/2504.07405](https://arxiv.org/abs/2504.07405).

564 Zhipeng Huang, Shaobin Zhuang, Canmiao Fu, Binxin Yang, Ying Zhang, Chong Sun, Zhizheng  
565 Zhang, Yali Wang, Chen Li, and Zheng-Jun Zha. Wegen: A unified model for interactive multi-  
566 modal generation as we chat, 2025b. URL <https://arxiv.org/abs/2503.01115>.

567

568 Minguk Kang, Jun-Yan Zhu, Richard Zhang, Jaesik Park, Eli Shechtman, Sylvain Paris, and Taesung  
569 Park. Scaling up gans for text-to-image synthesis. In *Proceedings of the IEEE/CVF conference on  
570 computer vision and pattern recognition*, pp. 10124–10134, 2023.

571

572 Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete  
573 Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings  
574 of the IEEE/CVF international conference on computer vision*, pp. 4015–4026, 2023.

575

576 Yuval Kirstain, Adam Polyak, Uriel Singer, Shahbuland Matiana, Joe Penna, and Omer Levy. Pick-  
577 a-pic: An open dataset of user preferences for text-to-image generation. *Advances in Neural  
578 Information Processing Systems*, 36:36652–36663, 2023.

579

580 Nupur Kumari, Bingliang Zhang, Richard Zhang, Eli Shechtman, and Jun-Yan Zhu. Multi-concept  
581 customization of text-to-image diffusion. In *Proceedings of the IEEE/CVF conference on computer  
582 vision and pattern recognition*, pp. 1931–1941, 2023.

583

584 Duong H. Le, Tuan Pham, Sangho Lee, Christopher Clark, Aniruddha Kembhavi, Stephan Mandt,  
585 Ranjay Krishna, and Jiasen Lu. One diffusion to generate them all, 2024. URL [https://  
586 arxiv.org/abs/2411.16318](https://arxiv.org/abs/2411.16318).

587

588 Dongxu Li, Junnan Li, and Steven Hoi. Blip-diffusion: Pre-trained subject representation for  
589 controllable text-to-image generation and editing. *Advances in Neural Information Processing  
590 Systems*, 36:30146–30166, 2023a.

591

592 Zhen Li, Mingdeng Cao, Xintao Wang, Zhongang Qi, Ming-Ming Cheng, and Ying Shan. Photomaker:  
593 Customizing realistic human photos via stacked id embedding, 2023b. URL [https://arxiv.  
594 org/abs/2312.04461](https://arxiv.org/abs/2312.04461).

595

596 Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr  
597 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer vision–  
598 ECCV 2014: 13th European conference, zurich, Switzerland, September 6-12, 2014, proceedings,  
599 part v 13*, pp. 740–755. Springer, 2014.

---

594 Zheyuan Liu, Cristian Rodriguez-Opazo, Damien Teney, and Stephen Gould. Image retrieval on  
595 real-life images with pre-trained vision-and-language models. In *Proceedings of the IEEE/CVF*  
596 *International Conference on Computer Vision*, pp. 2125–2134, 2021.

597

598 Zhiheng Liu, Yifei Zhang, Yujun Shen, Kecheng Zheng, Kai Zhu, Ruili Feng, Yu Liu, Deli Zhao,  
599 Jingren Zhou, and Yang Cao. Cones 2: Customizable image synthesis with multiple subjects, 2023.  
600 URL <https://arxiv.org/abs/2305.19327>.

601

602 Jian Ma, Junhao Liang, Chen Chen, and Haonan Lu. Subject-diffusion:open domain personalized  
603 text-to-image generation without test-time fine-tuning, 2024. URL <https://arxiv.org/abs/2307.11410>.

604

605 Zhendong Mao, Mengqi Huang, Fei Ding, Mingcong Liu, Qian He, and Yongdong Zhang. Re-  
606 alcustom++: Representing images as real-word for real-time customization. *arXiv preprint*  
607 *arXiv:2408.09744*, 2024.

608

609 OpenAI. Introducing gpt-4o and more tools to chatgpt free users, 2024. URL <https://openai.com/index/gpt-4o-and-more-tools-to-chatgpt-free/>. Accessed: 2024-06-15.

610

611 Or Patashnik, Rinon Gal, Daniil Ostashev, Sergey Tulyakov, Kfir Aberman, and Daniel Cohen-  
612 Or. Nested attention: Semantic-aware attention values for concept personalization, 2025. URL  
613 <https://arxiv.org/abs/2501.01407>.

614

615 Maitreya Patel, Sangmin Jung, Chitta Baral, and Yezhou Yang.  $\lambda$ -eclipse: Multi-concept personalized  
616 text-to-image diffusion models by leveraging clip latent space. *arXiv preprint arXiv:2402.05195*,  
617 2024.

618

619 William Peebles and Saining Xie. Scalable diffusion models with transformers, 2023. URL <https://arxiv.org/abs/2212.09748>.

620

621 Yuang Peng, Yuxin Cui, Haomiao Tang, Zekun Qi, Runpei Dong, Jing Bai, Chunrui Han, Zheng Ge,  
622 Xiangyu Zhang, and Shu-Tao Xia. Dreambench++: A human-aligned benchmark for personalized  
623 image generation. In *The Thirteenth International Conference on Learning Representations*, 2025.  
624 URL <https://openreview.net/forum?id=4GSOESJrk6>.

625

626 pin. <https://www.pinterest.com/>. <https://www.pinterest.com/>.

627

628 Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe  
629 Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image  
630 synthesis, 2023. URL <https://arxiv.org/abs/2307.01952>.

631

632 Zeju Qiu, Weiyang Liu, Haiwen Feng, Yuxuan Xue, Yao Feng, Zhen Liu, Dan Zhang, Adrian Weller,  
633 and Bernhard Schölkopf. Controlling text-to-image diffusion by orthogonal finetuning. *Advances*  
634 *in Neural Information Processing Systems*, 36:79320–79362, 2023.

635

636 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,  
637 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual  
638 models from natural language supervision. In *International conference on machine learning*, pp.  
639 8748–8763. PmLR, 2021.

640

641 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
642 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer-*  
643 *ence on computer vision and pattern recognition*, pp. 10684–10695, 2022a.

644

645 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
646 resolution image synthesis with latent diffusion models, 2022b. URL <https://arxiv.org/abs/2112.10752>.

647

648 Ciara Rowles, Shimon Vainer, Dante De Nigris, Slava Elizarov, Konstantin Kutsy, and Simon Donné.  
649 Ipadapter-instruct: Resolving ambiguity in image-based conditioning using instruct prompts, 2024.  
650 URL <https://arxiv.org/abs/2408.03209>.

---

648 Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman.  
649 Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *Proceed-*  
650 *ings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 22500–22510,  
651 2023.

652 Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar  
653 Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic  
654 text-to-image diffusion models with deep language understanding. *Advances in neural information*  
655 *processing systems*, 35:36479–36494, 2022.

656  
657 Jing Shi, Wei Xiong, Zhe Lin, and Hyun Joon Jung. Instantbooth: Personalized text-to-image genera-  
658 tion without test-time finetuning, 2023. URL <https://arxiv.org/abs/2304.03411>.

659  
660 Quan Sun, Yufeng Cui, Xiaosong Zhang, Fan Zhang, Qiyang Yu, Yueze Wang, Yongming Rao,  
661 Jingjing Liu, Tiejun Huang, and Xinlong Wang. Generative multimodal models are in-context  
662 learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,  
663 pp. 14398–14409, 2024.

664  
665 Shangkun Sun, Bowen Qu, Xiaoyu Liang, Songlin Fan, and Wei Gao. Ie-bench: Advancing  
666 the measurement of text-driven image editing for human perception alignment. *arXiv preprint*  
667 *arXiv:2501.09927*, 2025.

668  
669 Zhenxiong Tan, Songhua Liu, Xingyi Yang, Qiaochu Xue, and Xinchao Wang. Ominicontrol:  
670 Minimal and universal control for diffusion transformer. *arXiv preprint arXiv:2411.15098*, 2024.  
671 uns. <https://unsplash.com/>. <https://unsplash.com/>.

672  
673 Andrey Voynov, Qinghao Chu, Daniel Cohen-Or, and Kfir Aberman. P+: Extended textual condition-  
674 ing in text-to-image generation, 2023. URL <https://arxiv.org/abs/2303.09522>.

675  
676 Haofan Wang, Matteo Spinelli, Qixun Wang, Xu Bai, Zekui Qin, and Anthony Chen. Instantstyle:  
677 Free lunch towards style-preserving in text-to-image generation. *arXiv preprint arXiv:2404.02733*,  
678 2024a.

679  
680 Xierui Wang, Siming Fu, Qihan Huang, Wanggui He, and Hao Jiang. Ms-diffusion: Multi-subject  
681 zero-shot image personalization with layout guidance. *arXiv preprint arXiv:2406.07209*, 2024b.

682  
683 Yibin Wang, Yuhang Zang, Hao Li, Cheng Jin, and Jiaqi Wang. Unified reward model for multimodal  
684 understanding and generation. *arXiv preprint arXiv:2503.05236*, 2025.

685  
686 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi,  
687 Quoc V. Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language  
688 models. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.

689  
690 Yuxiang Wei, Yabo Zhang, Zhilong Ji, Jinfeng Bai, Lei Zhang, and Wangmeng Zuo. Elite: Encoding  
691 visual concepts into textual embeddings for customized text-to-image generation. In *Proceedings*  
692 *of the IEEE/CVF International Conference on Computer Vision*, pp. 15943–15953, 2023.

693  
694 Shaojin Wu, Mengqi Huang, Wenxu Wu, Yufeng Cheng, Fei Ding, and Qian He. Less-to-more general-  
695 ization: Unlocking more controllability by in-context generation. *arXiv preprint arXiv:2504.02160*,  
696 2025.

697  
698 Xiaoshi Wu, Keqiang Sun, Feng Zhu, Rui Zhao, and Hongsheng Li. Human preference score:  
699 Better aligning text-to-image models with human preference. In *Proceedings of the IEEE/CVF*  
700 *International Conference on Computer Vision*, pp. 2096–2105, 2023.

701  
702 Shitao Xiao, Yueze Wang, Junjie Zhou, Huaying Yuan, Xingrun Xing, Ruiran Yan, Chaofan Li,  
703 Shuting Wang, Tiejun Huang, and Zheng Liu. Omnigen: Unified image generation. *arXiv preprint*  
*arXiv:2409.11340*, 2024.

704  
705 Zhexiao Xiong, Wei Xiong, Jing Shi, He Zhang, Yizhi Song, and Nathan Jacobs. Grounding-  
706 booth: Grounding text-to-image customization, 2025. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2409.08520)  
[2409.08520](https://arxiv.org/abs/2409.08520).

---

702 Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao Dong.  
703 Imagereward: Learning and evaluating human preferences for text-to-image generation. *Advances*  
704 *in Neural Information Processing Systems*, 36:15903–15935, 2023.  
705

706 Jiazheng Xu, Yu Huang, Jiale Cheng, Yuanming Yang, Jiajun Xu, Yuan Wang, Wenbo Duan,  
707 Shen Yang, Qunlin Jin, Shurun Li, Jiayan Teng, Zhuoyi Yang, Wendi Zheng, Xiao Liu, Ming  
708 Ding, Xiaohan Zhang, Xiaotao Gu, Shiyu Huang, Minlie Huang, Jie Tang, and Yuxiao Dong.  
709 Visionreward: Fine-grained multi-dimensional human preference learning for image and video  
710 generation, 2024. URL <https://arxiv.org/abs/2412.21059>.

711 Hu Ye, Jun Zhang, Sibao Liu, Xiao Han, and Wei Yang. Ip-adapter: Text compatible image prompt  
712 adapter for text-to-image diffusion models. *arXiv preprint arXiv:2308.06721*, 2023.  
713

714 Yu Zeng, Vishal M. Patel, Haochen Wang, Xun Huang, Ting-Chun Wang, Ming-Yu Liu, and Yogesh  
715 Balaji. Jedi: Joint-image diffusion models for finetuning-free personalized text-to-image generation,  
716 2024. URL <https://arxiv.org/abs/2407.06187>.

717 Yuxuan Zhang, Yiren Song, Jiaming Liu, Rui Wang, Jinpeng Yu, Hao Tang, Huaxia Li, Xu Tang,  
718 Yao Hu, Han Pan, et al. Ssr-encoder: Encoding selective subject representation for subject-  
719 driven generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*  
720 *Recognition*, pp. 8069–8078, 2024.  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

## A DETAILS OF HIERARCHICAL CATEGORY ESTABLISHING

**The First-level Category** We observed the composition of existing benchmark data. From a more abstract and higher-level perspective, images in these datasets could be categorized into two types: photorealistic and non-photorealistic. Theoretically, the specific image categories represented within these two types can be identical. To maintain consistency with previous work and to ensure comprehensive data sampling, we designated photorealistic and non-photorealistic as the first-level categories. Furthermore, we ensure that the specific subcategories under both photorealistic and non-photorealistic types are fully aligned.

**The Second-level Category** We examined both the DreamBench and DreamBench++ datasets. In DreamBench, the dataset is divided into two categories: living subjects and objects. DreamBench++ further refines this categorization by introducing three categories: living subjects, objects, and style. We construct our secondary subcategories based on them. We define our secondary categories as objects, humans, and animals. Specifically, we subdivide the "living subjects" category into "humans" and "animals," as humans exhibit significantly different visual characteristics compared to animals. For the human category, we place particular emphasis on the accuracy of facial feature reconstruction, acknowledging the existence of dedicated research domains focused on facial preservation. In contrast, animals generally display greater variability in appearance than human faces. In comparison to DreamBench++, we exclude the "style" category. This decision is motivated by the focus of our task on subject-driven T2I generation, where "style" does not constitute a tangible entity. Moreover, including the style category would complicate the calculation of subject consistency, whereas our work is primarily concerned with the customization of entities.

**The Third-level Category** For the third-level categories, our objective was to strike a balance between granularity and generality. Categories that are too broad may result in insufficient keyword retrieval, potentially introducing bias into the final image sampling. Conversely, overly fine-grained categories may hinder subsequent experimental analysis by diluting meaningful insights. To address this, we consulted existing large-scale datasets such as COCO and ImageNet, as well as Wikipedia, to compile a list of candidate category labels. The specific labels are listed in Table 4. This comprehensive set of labels ensured broad coverage. However, many of these labels were excessively detailed, so we employ GPT-4o to merge them, followed by manual review to ensure the rationality and coherence of the final categories. The correspondence between the third-level categories and the candidate category labels is presented in Table 4. For the "human" category, we introduced a specific distinction by dividing it into "celebrities & artistic figures," "facial close-ups," and "half-body or full-body photo". We observed that models tend to perform significantly better on celebrities, which we hypothesize is due to the inclusion of celebrity data in the training sets of text-to-image foundation models. Table 14 provides empirical support for our hypothesis to some extent. The rationale for distinguishing between facial close-ups and non-facial close-ups is that the former focuses exclusively on the facial details of the individual in the reference image, whereas the latter also requires attention to the body details.

Through the aforementioned steps, we constructed a hierarchical category system. The resulting category hierarchy is presented in Figure 8.

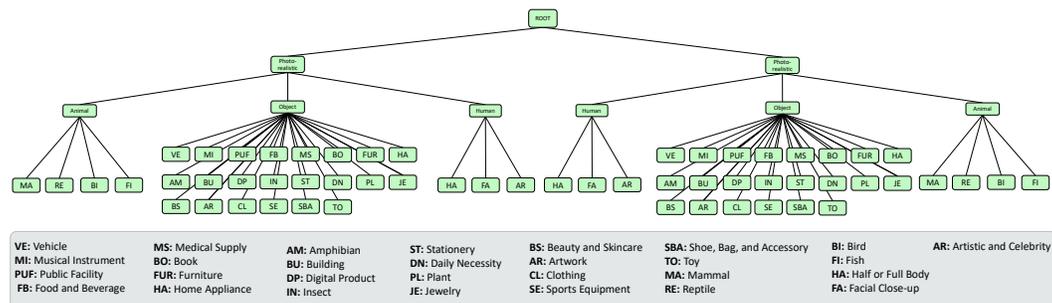


Figure 8: **The hierarchical category system.** We developed a three-level category hierarchy by integrating data from existing large-scale datasets and open-source encyclopedic resources.

810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863

Table 4: The correspondence between the third-level categories and the candidate category labels

Candidate Category Labels							The Third-level Category
reptile	lizard	dinosaur	turtle	crocodile	chameleon	gecko	Reptile
fly	firefly	ant	butterfly	ladybug	locust	dragonfly	Insect
amphibian	frog	bullfrog	toad	salamander			Amphibian
fish	goldfish	seahorse	shark	tilapia			Fish
bird	chicken	duck	owl	swan	goose	rooster	Bird
hen	turkey	swallow	crow	pigeon			
mammal	cat	dog	horse	sheep	cow	elephant	Mammal
bear	squirrel	giraffe	lion	monkey	tiger	bunny	
goat	pig	kangaroo	rhinoceros	deer	hippo	platypus	
whale	aardvark	rabbit	zebra	mouse			
street	fountain	fire hydrant	traffic light	sign	parking meter	goal net	Public Facility
field goal post	soccer net	basketball court	bus stop sign				
furniture	dining table	sofa	chair	couch	bed	desk	Furniture
table	coffee table	side table	bench	cabinet	mirror	carpet	
window	door	chandelier	table lamp	gate			
flower	potted plant	tree	sunflower	cactus	lavender		Plant
cookie	milk	pancake	pasta	grape	cereal	bean	Food and Beverage
pineapple	carrot	broccoli	banana	orange	strawberry	apple	
bread	sandwich	cake	pizza	soup	meat	pumpkin	
cheese	cupcake	donut	hot dog	bacon	egg	tomato	
dryer	fridge	refrigerator	microwave	oven	toaster	washer	Home Appliance
blender	hair drier	fan (ceiling/floor)	printer	fax machine	copier		
necklace	bracelet	ring	pendant	brooch	anklet		Jewelry
wheelchair	gauze	crutch	stethoscope	syringe			Medical Supply
pants	jacket	long sleeve shirt	short sleeve shirt	pajamas	underpants	shirt	Clothing
shorts	scarf	tie	super hero costume	sock			
book	magazine	textbook	dictionary	biography			Book
bat	skis	snowboard	tennis racket	basketball hoop	baseball glove	soccer ball	Sports Equipment
sports ball	basketball	football	tennis net	hoop			
flip flop	handbag	glove	shoe	backpack			Shoe, Bag, and Accessory
pen	pencil	fax machine	stapler				Stationery
vehicle	car	van	truck	bus	train	boat	Vehicle
sailboat	raft	airplane	helicopter	hot air balloon	rocket	bicycle	
unicycle	motorcycle	motorbike	skateboard				
house	building	roof	bridge	church			Building
picture frame	movie (disc)	playing cards	table cloth				Artwork
musical instrument	guitar	drum	flute	violin			Musical Instrument
telephone	laptop	computer	tablet	ipad	iphone	cell phone	Digital Product
remote	mouse	keyboard	printer	desktop	copier	radio	
kite	toy cars	toy	legos	robot	doll		
hair brush	toner	blush	serum	emulsion	sunscreen		Beauty and Skincare
bottle	plate	cup	bowl	teapot	fork	knife	Daily Necessity
spoon	clock	toothbrush	vase	towel	candle	balloon	
box	chopping board	ladder	basket	pillow	power outlet	light switch	
person							Person

## B DETAILS OF KEYWORDS COLLECTION

The keywords utilized during the image collection process are presented in Table 5. During the keyword collection process, we utilized the following prompt for GPT-4o:

*"You are a researcher with extensive knowledge of various real-world entity classifications. Given a specific category, please generate detailed, non-redundant instances relevant to this category. The category is {}.*  
*The corresponding instances are as follows:"*

864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917

Table 5: Based on the categories, we employ GPT-4o to generate keyword associations and further enhanced the results by incorporating manually curated keywords.

The Third-level Category	Keywords						
Vehicle	van pickup truck	steam locomotive bicycle	car boat	airplane taxi	UFO motorcycle	hot air balloon subway	oil tanker
Musical Instrument	guitar pick suona	electronic drum saxophone	digital piano harmonica	guitar cello	snare drum violin	flute pipa	african drum erhu
Public Facility	fire extinguisher	traffic sign	street lamp	street	station		
Food and Beverage	edible oil pineapple apple vegetable	instant noodles milk donut chicken	water orange durian noodles	pastries avocado sports drink hamburger	coffee can canned health products salad	biscuits juice egg chocolate	edible salt milk powder rice yogurt
Medical Supply	band-aid blood glucose meter	medicine crutch	wheelchair stethoscope	disinfectant syringe	first aid kit	medication	medicine bottle
Book	yearbook book	almanac notebook	workbook magazine	comic dictionary	encyclopedia	atlas	pamphlet
Furniture	shelf barber chair ottoman	makeup mirror office chair bookcase	stool bathroom mirror wardrobe	bathroom cabinet chair nightstand	cabinet sofa dresser	bean bag chair dining table	children's chair bed
Home Appliance	beauty device microphone television	kettle refrigerator oven	speaker hair dryer juicer	massage chair humidifier dishwasher	vacuum cleaner washing machine	rice cooker microwave oven	robot vacuum curving iron
Amphibian	newt frog	olm toad	bullfrog caecilian	wood frog salamander	Surinam toad	alpine newt	glass frog
Building	house hut	apartment building leaning tower of pisa	duplex house pyramid	church statue of liberty	temple of heaven eiffel tower	castle	golden gate bridge
Digital Product	smart robot printer smartwatch	headphones camcorder vintage camera	e-book reader camera monitor	desktop computer smart camera drone	roll of film laptop projector	router mobile phone fitness tracker	tablet walkie-talkie
Insect	shrimp	crab	ant	grasshopper	butterfly		
Stationery	glue stick stapler	globe crayon	calculator ballpoint pen	floppy disk eraser	tape measure	scissors	compass
Daily Necessity	hammer birdcage glass jar electric saw	candle alarm clock vase mop	mug spoon hanger broom	teapot bowl soap dish comb	berry bowl toothbrush frying pan	curtain shower gel baby bottle	pillow clock kitchen knife
Plant	cactus mint	coconut tree rose	tree sunflower	potted plant tulip	peony cactus	willow tree lavender	maple leaf
Jewelry	earrings tiara gold bar	ring crown necklace	crystal stud pendant	bracelet chain brooch	watch gemstone anklet	hair accessory choker locket	beaded bracelet hairpin
Beauty and Skincare	perfume blush	makeup brush eye shadow	lotion facial serum	sunscreen spray emulsion	face cream serum	nail polish mascara	toner lipstick
Artwork	bouquet of flowers sculpture	clay sculpture ceramic craft	wood carving mural	classical bust relief	stone carving	catstatue	mugskulls
Clothing	dress pants	baby clothes shirt	clothing down jacket	jeans coat	sweatshirt skirt	T-shirt shorts	socks vest
Sports Equipment	tennis adjustable bench treadmill	ball knee pad skateboard	tent backpack barbell	trekking poles soccer dumbbell	yoga mat sleeping bag	billiard baseball	badminton flamingo float
Shoe, Bag and Accessory	suitcase glasses hat	slippers sandals backpack	sunglasses shoes cap	canvas shoes luggage purse tie	high-top shoes fancy boot handbag	sports shoes belt sandals	scarf sneaker
Toy	actionfigure robot minion	monster toy motorbike toy smart robot	car magic cube robot toy	egg poop emoji toy	duck toy sloth plushie wolf plushie	teddy bear bear plushie doll	balloon red cartoon Eevee figurine
Mammal	rabbit panda alpaca	fox elephant puppy	wolf llama monkey	Siamese cat tiger kitten	polar bear dog dolphin	cat raccoon French bulldog	deer lion
Reptile	cobra turtle	gecko sea turtle	rattlesnake soft-shelled turtle	crocodile snake	chameleon lizard	alligator	iguana
Bird	heron woodpecker peacock bird	pigeon nightingale swallow canary	toucan duck owl sparrow	parrot turkey kingfisher rooster	stork chicken hawk	flamingo crow dove	penguin eagle anchovy
Fish	shark skate	tropical fish swordfish	jellyfish herring	goldfish sardine	perch carp	eel salmon	monkfish tuna
Person	person						

## C DETAILS OF PROMPT GENERATION

The specific instructions used in prompt generation are detailed in Figure 5. During the actual generation process, some of the prompts produced by GPT-4o did not meet the required criteria. Therefore, we instructed GPT-4o to generate multiple prompts for each image, and then manually selected those that best matched the intended scenarios. Figure 14 presents the results generated by different methods in this study, along with their corresponding prompts.

---

## 918 D ADDITIONAL DISCUSSIONS AND DETAILS OF MODEL PERFORMANCE

### 919 D.1 ADDITIONAL DISCUSSIONS

920  
921  
922 **Analysis of The First-Level Category** The primary categories are divided into photorealistic and  
923 non-photorealistic. Table 6 and Figure 9 present the performance of different methods on these  
924 two categories across three evaluation dimensions. The results show that: (1) *Subject Preservation*:  
925 Almost all methods perform better on photorealistic categories than on non-photorealistic ones.  
926 We speculate that this is because, when referencing subjects from non-photorealistic categories,  
927 these methods tend to generate photorealistic images based on the prompt, which results in lower  
928 subject consistency. (2) *Prompt Following*: The performance gap between photorealistic and non-  
929 photorealistic categories is relatively small. This can be attributed to the fact that CLIP-T focuses  
930 primarily on the semantic information of the image. As long as the generated subject matches the  
931 category and general appearance described in the prompt, the CLIP-T score will not be significantly  
932 reduced. (3) *Image quality*: There is little difference in performance between photorealistic and  
933 non-photorealistic categories. This indicates that the distinction between these two categories does  
934 not affect the quality of image generation, and the HPSv2 metric does not show a preference for  
935 either category.

936 **Analysis of The Second-Level Category** The secondary categories under both the realistic and  
937 non-realistic primary categories are further subdivided into objects, humans, and animals. Table 7  
938 and Figure 10 present the performance of various methods across these three dimensions for both  
939 realistic and non-realistic categories. The results demonstrate that, irrespective of whether the primary  
940 category is realistic or non-realistic, the scores for the subject preservation dimension are consistently  
941 lower for the human category across nearly all models. As detailed in Table 8, this phenomenon can  
942 be attributed to the distribution of difficulty levels within the human category, where the proportions  
943 of simple, medium, and hard cases are 1.96%, 50.98%, and 47.06%, respectively. In contrast, the  
944 object and animal categories exhibit a higher proportion of subjects at the simple difficulty level and a  
945 lower proportion at the hard difficulty level, which likely contributes to their relatively higher subject  
946 preservation scores.

947 **Implications for Technical Approaches** (1) Figure 11 shows that, as base models and model  
948 architectures are updated, the performance boundary of these models consistently expands outward.  
949 Table 9 presents all the base models used by each method. It can be observed that the top-performing  
950 methods consistently employ relatively recent text-to-image base models. For instance, UNO utilizes  
951 FLUX as its foundational model. This observation suggests that the adoption of advanced text-to-  
952 image base models is a critical factor in enhancing performance on subject-driven T2I tasks. (2)  
953 Historically, fine-tuning methods have generally outperformed encoder-based approaches in terms  
954 of subject preservation. This advantage is attributed to their ability to better retain the original  
955 text-image conditional distribution by fine-tuning on images of the specified subject. In contrast,  
956 encoder-based methods often encounter interference during feature injection, which can hinder precise  
957 prompt alignment. However, with the development of more advanced encoding techniques, the  
958 adoption of larger and more powerful base models, and the availability of extensive training datasets,  
959 encoder-based methods have demonstrated significantly improved performance. From an application  
960 standpoint, fine-tuning methods require substantial computational resources for optimization and often  
961 exhibit limited generalization capabilities. In contrast, encoder-based methods are less constrained  
962 by these limitations, making them more practical for future applications. Nevertheless, our analysis  
963 indicates that current encoder-based methods still face challenges in accurately reconstructing subjects  
964 with high-frequency details in images. This limitation may stem from the characteristics of commonly  
965 used image encoders, such as CLIP, which tend to prioritize semantic information over fine-grained  
966 details. Consequently, future research should focus on enhancing the restoration of challenging  
967 subject details.

### 967 D.2 DETAILS OF MODEL PERFORMANCE

968  
969 In this section, we present the detailed evaluation results for each metric across all models. To  
970 comprehensively evaluate the effectiveness of different metrics for assessing subject consistency, we  
971 calculated multiple metrics for each method. The detailed results are presented in Table 10, 11, 12.  
In section 4, we present the performance of all methods across images with different difficulty

972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025

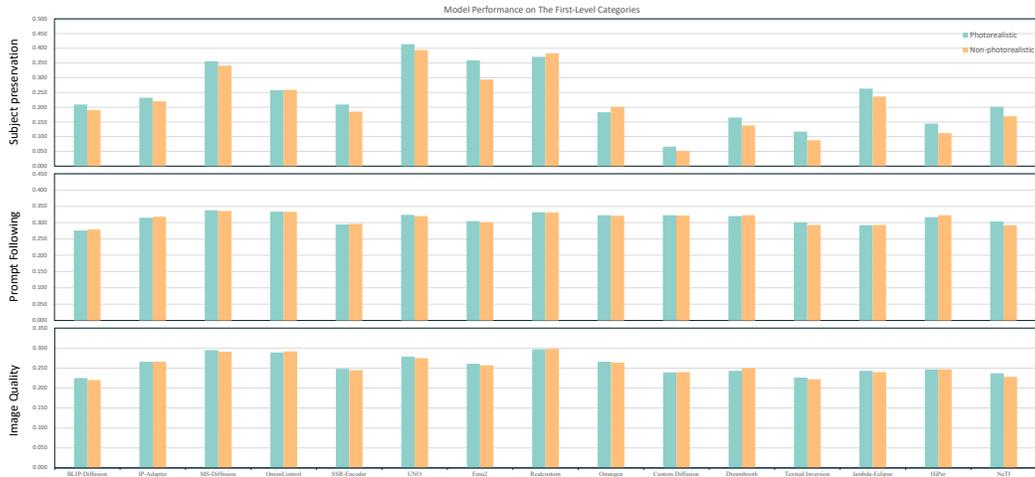


Figure 9: Comparison of bar charts for DSH-Bench scores in different first-level categories.

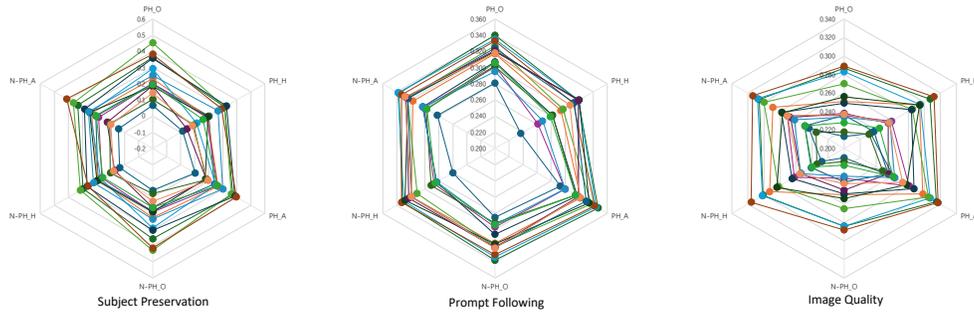


Figure 10: Comparison of radar charts for DSH-Bench scores in different second-level categories.

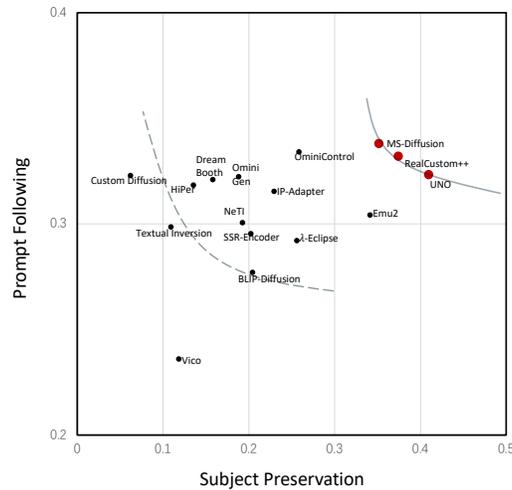


Figure 11: Pareto front diagram illustrating model performance across both subject and prompt dimensions. The red points in the diagram represent the current Pareto-optimal solutions.

levels, different prompt scenarios, and multiple categories. We show the specific metric values in Table 13, 14, 15, 16. Table 9 shows the full ranking among all methods.

Table 6: We evaluate the performance of various methods on DSH-Bench dataset, specifically analyzing their effectiveness across **the first-level categories**. PH: Photorealistic. N-PH: Non-Photorealistic.

Method	Subject Preservation $\uparrow$		Prompt Following $\uparrow$		Image Quality $\uparrow$	
	PH	N-PH	PH	N-PH	PH	N-PH
BLIP-Diffusion	0.209	0.190	0.276	0.279	0.225	0.220
IP-Adapter	0.232	0.220	0.315	0.318	0.266	0.266
MS-Diffusion	0.356	0.341	<b>0.338</b>	<b>0.336</b>	0.295	0.291
OminiControl	0.258	0.259	0.334	0.333	0.289	0.292
SSR-Encoder	0.209	0.185	0.295	0.296	0.248	0.245
UNO	<b>0.414</b>	<b>0.394</b>	0.324	0.320	0.279	0.275
Emu2	0.359	0.294	0.305	0.301	0.261	0.257
RealCustom++	0.371	0.383	0.332	0.331	<b>0.297</b>	<b>0.298</b>
OmniGen	0.183	0.201	0.323	0.321	0.266	0.264
Custom Diffusion	0.066	0.052	0.323	0.322	0.239	0.240
DreamBooth	0.165	0.138	0.320	0.323	0.243	0.250
Textual Inversion	0.117	0.088	0.301	0.293	0.226	0.222
$\lambda$ -Eclipse	0.263	0.236	0.292	0.293	0.243	0.240
HiPer	0.144	0.112	0.317	0.323	0.247	0.247
NeTI	0.201	0.169	0.304	0.292	0.237	0.228
<i>Aver.</i>	0.236	0.217	0.313	0.312	0.257	0.256

## E IMPLEMENTATION DETAILS

### E.1 EXPERIMENTAL DETAILS OF EXISTING METHODS

The configurations for the training hyperparameters used in training-based methods on DSH-Bench are detailed in Table 17. To ensure a fair comparison in inference stage, we generated four images for each prompt of every image. The final evaluation metrics were calculated as the average score across these four images.

### E.2 DETAILS OF SICS IMPLEMENTATION

**How the SICS metric is computed** Rather than relying on simple embedding distances, SICS instruction-tunes multimodal large language models to directly produce fine-grained subject-consistency scores (0–5) with accompanying explanations. These criteria align closely with human judgments. As illustrated in Figure 12, SICS employs prompts during fine-tuning that explicitly target subject consistency while de-emphasizing global image semantics. This design reduces background and style confounds that bias CLIP-like methods, yielding tighter alignment with the core requirements of subject-consistency evaluation. Consequently, SICS focuses on core visual attributes rather than high-level semantics. A further advantage of SICS is its scoring granularity, which mitigates the “score saturation” phenomenon commonly observed in the upper range of GPT-4o evaluations. Figure 2 shows representative cases in which SICS aligns more closely with human assessments.

**Evaluation Instruction** Figure 12 illustrates the annotation criteria of the training dataset as well as the training process.

**Datasets** We collected a substantial number of image pairs. To ensure data quality, we applied standardized filtering and preprocessing procedures, such as enforcing a minimum image resolution of 512 pixels. Additionally, we employed Qwen2.5-VL-72B to conduct preliminary screening. After this automated filtering, five annotators manually annotated the remaining image pairs according to the guidelines illustrated in Figure 12.

**Training Details** We fine-tuned Qwen2.5-VL-7B on the manually annotated dataset described above. All experiments were conducted using 8 GPUs. For the learning rate, we experimented with the set  $1e5$ . The batch size per device was set to 4, with a gradient accumulation step of 8.

1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133

Table 7: We evaluate the performance of various methods on DSH-Bench dataset, specifically analyzing their effectiveness across **the second-level categories**. PH: Photorealistic, N-PH: Non-Photorealistic, O: Object, A: Animal, H: Human.

Method	Subject Preservation					
	PH_O	PH_H	PH_A	N-PH_O	N-PH_H	N-PH_A
BLIP-Diffusion	0.202	0.201	0.24	0.186	0.189	0.206
IP-Adapter	0.232	0.193	0.267	0.226	0.188	0.237
MS-Diffusion	0.362	0.315	0.371	0.358	0.296	0.333
OminiControl	0.293	0.114	0.249	0.291	0.17	0.247
SSR-Encoder	0.199	0.186	0.26	0.193	0.162	0.185
UNO	<b>0.453</b>	0.312	0.361	<b>0.428</b>	<b>0.315</b>	0.365
Emu2	0.358	<b>0.326</b>	0.387	0.305	0.266	0.285
RealCustom++	0.383	0.291	<b>0.396</b>	0.415	0.26	<b>0.412</b>
OmniGen	0.183	0.194	0.176	0.19	0.196	0.249
Custom Diffusion	0.067	0.014	0.103	0.059	0.035	0.043
DreamBooth	0.188	0.044	0.184	0.164	0.07	0.124
Textual Inversion	0.104	0.091	0.184	0.078	0.101	0.105
$\lambda$ -Eclipse	0.252	0.266	0.3	0.236	0.221	0.256
HiPer	0.143	0.083	0.195	0.126	0.079	0.098
NeTI	0.195	0.159	0.259	0.164	0.156	0.201
<i>Aver.</i>	0.241	0.186	0.262	0.228	0.180	0.223
Method	Prompt Following					
	PH_O	PH_H	PH_A	N-PH_O	N-PH_H	N-PH_A
BLIP-Diffusion	0.281	0.237	0.293	0.285	0.26	0.282
IP-Adapter	0.317	0.294	0.322	0.317	0.319	0.317
MS-Diffusion	<b>0.340</b>	0.319	<b>0.347</b>	<b>0.338</b>	0.332	0.337
OminiControl	0.335	0.319	0.344	0.334	0.33	<b>0.338</b>
SSR-Encoder	0.302	0.261	0.3	0.297	0.287	0.301
UNO	0.327	0.297	0.337	0.321	0.311	0.325
Emu2	0.307	0.282	0.317	0.306	0.283	0.303
RealCustom++	0.333	0.312	0.342	0.331	<b>0.333</b>	0.333
OmniGen	0.320	<b>0.320</b>	0.334	0.318	0.328	0.324
Custom Diffusion	0.324	0.313	0.33	0.322	0.319	0.324
DreamBooth	0.321	0.319	0.319	0.322	0.323	0.327
Textual Inversion	0.301	0.282	0.315	0.292	0.291	0.298
$\lambda$ -Eclipse	0.295	0.268	0.3	0.294	0.283	0.303
HiPer	0.318	0.307	0.32	0.323	0.319	0.328
NeTI	0.306	0.279	0.315	0.294	0.285	0.297
<i>Aver.</i>	0.315	0.294	0.322	0.313	0.307	0.316
Method	Image Quality					
	PH_O	PH_H	PH_A	N-PH_O	N-PH_H	N-PH_A
BLIP-Diffusion	0.213	0.233	0.262	0.21	0.228	0.244
IP-Adapter	0.251	0.294	0.298	0.25	0.293	0.289
MS-Diffusion	0.287	0.307	0.315	0.284	0.301	0.306
OminiControl	0.283	0.295	0.307	0.284	0.302	0.308
SSR-Encoder	0.236	0.259	0.281	0.232	0.262	0.271
UNO	0.270	0.285	0.305	0.265	0.282	0.3
Emu2	0.249	0.284	0.287	0.249	0.265	0.278
RealCustom++	<b>0.289</b>	<b>0.312</b>	<b>0.317</b>	<b>0.288</b>	<b>0.316</b>	<b>0.314</b>
OmniGen	0.256	0.294	0.278	0.254	0.284	0.277
Custom Diffusion	0.236	0.237	0.255	0.236	0.241	0.249
DreamBooth	0.238	0.255	0.255	0.245	0.254	0.267
Textual Inversion	0.218	0.231	0.248	0.214	0.234	0.235
$\lambda$ -Eclipse	0.234	0.257	0.263	0.23	0.254	0.262
HiPer	0.237	0.256	0.273	0.238	0.255	0.271
NeTI	0.228	0.244	0.261	0.218	0.24	0.249
<i>Aver.</i>	0.248	0.270	0.280	0.246	0.267	0.275

1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187

Table 8: Subject hard level distribution under the second category

Benchmark	Photorealistic								
	Object			Human			Animal		
	Easy	Medium	Hard	Easy	Medium	Hard	Easy	Medium	Hard
DreamBench	3	10	7	0	0	0	0	7	2
DreamBench++	6	24	31	0	7	5	0	26	16
DSH-Bench	54	85	84	1	26	24	2	39	21

Benchmark	Non-photorealistic								
	Object			Human			Animal		
	Easy	Medium	Hard	Easy	Medium	Hard	Easy	Medium	Hard
DreamBench	1	0	0	0	0	0	0	0	0
DreamBench++	2	1	1	0	1	7	0	1	2
DSH-Bench	28	32	15	1	4	20	3	11	8

Table 9: The full DSH-Bench leaderboard. The models are ranked by the final score  $S_h$ .

Method	T2I Model	Subject Preservation	Prompt Following	Image Quality	$S_h \uparrow$
RealCustom++	SDXL	0.375	0.332	<b>0.298</b>	<b>0.110</b>
UNO	FLUX.1-dev	<b>0.409</b>	0.323	0.278	0.109
MS-Diffusion	SDXL	0.352	<b>0.338</b>	0.294	0.107
Emu2	SDXL	0.341	0.304	0.260	0.089
OminiControl	FLUX.1-schnell	0.258	0.334	0.290	0.085
IP-Adapter	SDXL	0.229	0.315	0.266	0.071
$\lambda$ -Eclipse	SDXL	0.256	0.292	0.242	0.069
OmniGen	SDXL	0.188	0.322	0.265	0.062
SSR-Encoder	SD v1.5	0.202	0.295	0.247	0.059
NeTI	SD v1.4	0.192	0.301	0.234	0.056
BLIP-Diffusion	SD v1.5	0.204	0.277	0.223	0.054
DreamBooth	SD v1.5	0.158	0.321	0.245	0.052
HiPer	SD v1.4	0.135	0.318	0.247	0.045
Textual Inversion	SD v1.5	0.109	0.299	0.225	0.035
ViCo	SD v1.4	0.118	0.236	0.186	0.029
Custom Diffusion	SD v1.4	0.062	0.323	0.240	0.023

Table 10: Evaluation of Subject-driven T2I generation model on DreamBench. C, D, Img, T and I represent CLIP, DINO, Image, Text and Image, respectively.

Method	Subject Preservation					Prompt Following		Image Quality		
	C-B-I $\uparrow$	C-L-I $\uparrow$	D-I $\uparrow$	D-v2-I $\uparrow$	SICS $\uparrow$	C-B-T $\uparrow$	C-L-T $\uparrow$	ImageReward $\uparrow$	PickScore $\uparrow$	HPSv2 $\uparrow$
BLIP-Diffusion	0.824	0.784	0.684	0.640	0.229	0.291	0.239	0.420	0.599	0.267
IP-Adapter	0.836	<b>0.820</b>	0.684	0.648	0.230	0.321	0.263	0.616	0.600	0.291
MS-Diffusion	0.814	0.796	0.732	0.687	0.316	0.332	0.279	0.775	0.600	0.311
OminiControl	0.784	0.772	0.614	0.555	0.279	<b>0.336</b>	<b>0.284</b>	0.793	0.593	0.306
SSR-Encoder	0.830	0.802	0.732	0.677	0.231	0.302	0.251	0.535	0.600	0.282
UNO	0.827	0.801	0.744	<b>0.716</b>	<b>0.409</b>	0.317	0.259	0.725	<b>0.602</b>	0.304
Emu2	<b>0.838</b>	0.818	0.737	0.704	0.360	0.291	0.235	0.463	0.599	0.272
RealCustom++	0.794	0.770	<b>0.746</b>	0.698	0.377	0.325	0.278	<b>0.813</b>	0.601	<b>0.316</b>

Table 11: Evaluation of Subject-driven T2I generation model on DreamBench++. C, D, Img, T and I represent CLIP, DINO, Image, Text and Image, respectively.

Method	Subject Preservation					Prompt Following		Image Quality		
	C-B-I $\uparrow$	C-L-I $\uparrow$	D-I $\uparrow$	D-v2-I $\uparrow$	SICS $\uparrow$	C-B-T $\uparrow$	C-L-T $\uparrow$	ImageReward $\uparrow$	PickScore $\uparrow$	HPSv2 $\uparrow$
BLIP-Diffusion	0.836	0.809	0.691	0.664	0.216	0.279	0.225	0.260	0.591	0.249
IP-Adapter	<b>0.846</b>	<b>0.845</b>	0.659	0.646	0.244	0.320	0.266	0.554	0.593	0.291
MS-Diffusion	0.812	0.823	0.666	0.653	0.346	<b>0.339</b>	<b>0.285</b>	0.729	0.593	0.309
OminiControl	0.761	0.780	0.551	0.566	0.268	0.336	0.284	<b>0.793</b>	0.593	0.308
SSR-Encoder	0.814	0.815	0.639	0.611	0.202	0.302	0.252	0.455	0.591	0.276
UNO	0.828	0.835	0.694	0.694	<b>0.410</b>	0.321	0.263	0.673	0.592	0.293
Emu2	0.833	0.823	0.665	0.632	0.343	0.309	0.255	0.460	0.593	0.275
RealCustom++	0.819	0.810	<b>0.714</b>	<b>0.706</b>	0.380	0.330	0.280	0.710	<b>0.594</b>	<b>0.314</b>

1188  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241

Table 12: Evaluation of Subject-driven T2I generation model on **DSH\_Bench**. C, D, Img, T and I represent CLIP, DINO, Image, Text and Image, respectively.

Method	Subject Preservation					Prompt Following			Image Quality		
	C-B-I↑	C-L-I↑	D-I↑	D-v2-I↑	SICS↑	C-B-T↑	C-L-T↑	ImageReward↑	PickScore↑	HPSv2↑	
BLIP-Diffusion	0.806	0.770	0.632	0.573	0.204	0.277	0.225	0.239	0.591	0.223	
IP-Adapter	0.824	0.812	0.610	0.577	0.229	0.315	0.263	0.493	0.594	0.266	
MS-Diffusion	0.786	0.783	0.623	0.600	0.352	<b>0.338</b>	0.287	0.705	<b>0.595</b>	0.294	
OminiControl	0.721	0.736	0.462	0.461	0.258	0.334	<b>0.288</b>	<b>0.787</b>	0.594	0.290	
SSR-Encoder	0.803	0.787	0.613	0.554	0.202	0.295	0.246	0.369	0.593	0.247	
UNO	0.781	0.784	0.607	0.599	<b>0.409</b>	0.323	0.272	0.705	0.594	0.278	
Emu2	0.815	0.804	0.631	0.606	0.341	0.304	0.256	0.441	0.594	0.260	
RealCustom++	0.781	0.769	0.645	0.624	0.374	0.332	0.285	0.695	0.595	<b>0.298</b>	
OminiGen	0.696	0.678	0.436	0.326	0.188	0.322	0.274	0.586	0.592	<b>0.265</b>	
Custom Diffusion	0.648	0.648	0.283	0.230	0.062	0.323	0.282	0.481	0.590	0.239	
DreamBooth	0.714	0.713	0.451	0.420	0.158	0.321	0.279	0.489	0.591	0.245	
Textual Inversion	0.689	0.683	0.372	0.320	0.109	0.299	0.253	0.340	0.590	0.225	
$\lambda$ -Eclipse	0.852	<b>0.833</b>	<b>0.676</b>	<b>0.638</b>	0.256	0.292	0.239	0.349	0.594	0.242	
HiPer	0.749	0.734	0.449	0.431	0.135	0.318	0.274	0.410	0.592	0.247	
NeTI	0.762	0.743	0.525	0.491	0.192	0.301	0.256	0.338	0.592	0.234	

Table 13: We evaluated the performance of various methods on DSH-Bench dataset, specifically analyzing their effectiveness across **the third-level categories (under photorealistic)**. Subject preservation, prompt following, and image quality are evaluated using SICS, CLIP-T, and HPSv2, respectively. **VE**: Vehicle, **MI**: Musical Instrument, **PUF**: Public Facility, **FB**: Food and Beverage, **MS**: Medical Supply, **BO**: Book, **FUR**: Furniture, **HA**: Home Appliance, **AM**: Amphibian, **BU**: Building, **DP**: Digital Product, **IN**: Insect, **ST**: Stationery, **DN**: Daily Necessity, **PL**: Plant, **JE**: Jewelry, **BS**: Beauty and Skincare, **AR**: Artwork, **CL**: Clothing, **SE**: Sports Equipment, **SBA**: Shoe, Bag, and Accessory, **TO**: Toy, **MA**: Mammal, **RE**: Reptile, **BI**: Bird, **FI**: Fish, **HF**: Half or Full Body, **FA**: Facial Close-up, **AC**: Artistic and Celebrity.

Method	Subject Preservation																												
	VE	MI	PUF	FB	MS	BO	FUR	HA	AM	BU	DP	IN	ST	DN	PL	JE	BS	AR	CL	SE	SBA	TO	MA	RE	BI	FI	HF	FA	AC
BLIP-Diffusion	0.246	0.181	0.142	0.182	0.151	0.187	0.220	0.195	0.242	0.285	0.192	0.185	0.157	0.212	0.257	0.164	0.176	0.189	0.247	0.225	0.209	0.205	0.181	0.202	0.202	0.194	0.219	0.202	0.194
IP-Adapter	0.217	0.244	0.183	0.217	0.207	0.137	0.287	0.228	0.208	0.300	0.216	0.123	0.257	0.230	0.248	0.189	0.232	0.162	0.292	0.246	0.262	0.218	0.303	0.190	0.248	0.225	0.195	0.221	0.139
MS-Diffusion	0.293	0.368	0.292	0.313	0.378	0.290	0.451	0.361	0.592	0.327	0.346	0.244	0.353	0.353	0.397	0.386	0.307	0.287	0.449	0.396	0.453	0.336	0.374	0.345	0.415	0.337	0.337	0.302	0.273
OminiControl	0.292	0.329	0.333	0.314	0.172	0.292	0.335	0.275	0.252	0.296	0.225	0.296	0.311	0.310	0.277	0.296	0.283	0.258	0.297	0.336	0.324	0.246	0.208	0.308	0.208	0.135	0.082	0.107	0.082
SSR-Encoder	0.187	0.237	0.156	0.191	0.201	0.150	0.315	0.202	0.192	0.286	0.183	0.146	0.158	0.195	0.293	0.112	0.158	0.143	0.204	0.231	0.185	0.203	0.301	0.197	0.209	0.215	0.159	0.228	0.188
UNO	0.387	0.411	0.446	0.444	0.461	0.330	0.580	0.442	0.442	0.461	0.442	0.325	0.402	0.457	0.427	0.396	0.418	0.419	0.515	0.502	0.524	0.436	0.359	0.322	0.420	0.300	0.367	0.258	0.247
Emu2	0.334	0.281	0.546	0.326	0.382	0.302	0.434	0.394	0.581	0.379	0.378	0.344	0.318	0.319	0.355	0.250	0.418	0.348	0.365	0.385	0.299	0.300	0.370	0.310	0.443	0.510	0.308	0.311	0.404
RealCustom++	0.304	0.421	0.371	0.340	0.328	0.228	0.479	0.395	0.500	0.406	0.347	0.327	0.358	0.399	0.432	0.349	0.325	0.347	0.389	0.488	0.369	0.417	0.394	0.342	0.450	0.390	0.306	0.270	0.283
OminiGen	0.111	0.123	0.094	0.249	0.161	0.183	0.196	0.133	0.058	0.236	0.192	0.096	0.157	0.169	0.229	0.146	0.251	0.179	0.202	0.173	0.203	0.180	0.197	0.082	0.203	0.144	0.209	0.215	0.113
Custom Diffusion	0.075	0.083	0.060	0.086	0.058	0.037	0.058	0.051	0.083	0.181	0.048	0.063	0.063	0.074	0.097	0.030	0.082	0.038	0.074	0.079	0.054	0.055	0.108	0.107	0.115	0.060	0.014	0.015	0.009
DreamBooth	0.225	0.180	0.206	0.206	0.151	0.077	0.224	0.232	0.200	0.242	0.158	0.190	0.195	0.220	0.307	0.135	0.163	0.105	0.131	0.187	0.167	0.186	0.181	0.170	0.190	0.208	0.048	0.035	0.050
Textual Inversion	0.143	0.081	0.090	0.116	0.081	0.022	0.099	0.085	0.192	0.225	0.099	0.117	0.115	0.089	0.174	0.115	0.074	0.101	0.071	0.114	0.087	0.134	0.191	0.133	0.219	0.160	0.082	0.078	0.138
$\lambda$ -Eclipse	0.280	0.215	0.117	0.215	0.214	0.170	0.346	0.248	0.283	0.312	0.253	0.156	0.211	0.249	0.275	0.243	0.211	0.193	0.298	0.283	0.309	0.238	0.340	0.227	0.272	0.231	0.261	0.273	0.268
HiPer	0.148	0.150	0.165	0.162	0.108	0.097	0.135	0.175	0.183	0.261	0.132	0.144	0.132	0.144	0.183	0.112	0.119	0.115	0.119	0.153	0.107	0.167	0.201	0.160	0.235	0.140	0.092	0.079	0.065
NeTI	0.199	0.154	0.140	0.195	0.179	0.123	0.235	0.200	0.292	0.275	0.144	0.196	0.182	0.225	0.264	0.176	0.151	0.188	0.145	0.193	0.193	0.226	0.262	0.252	0.282	0.233	0.152	0.139	0.218

Method	Prompt Following																													
	VE	MI	PUF	FB	MS	BO	FUR	HA	AM	BU	DP	IN	ST	DN	PL	JE	BS	AR	CL	SE	SBA	TO	MA	RE	BI	FI	HF	FA	AC	
BLIP-Diffusion	0.271	0.283	0.279	0.282	0.275	0.244	0.269	0.286	0.307	0.285	0.273	0.299	0.283	0.286	0.294	0.280	0.273	0.272	0.281	0.296	0.281	0.294	0.291	0.287	0.301	0.287	0.241	0.223	0.245	
IP-Adapter	0.315	0.332	0.310	0.318	0.314	0.301	0.306	0.324	0.331	0.312	0.320	0.319	0.320	0.317	0.316	0.306	0.307	0.310	0.319	0.325	0.321	0.333	0.325	0.303	0.327	0.306	0.291	0.291	0.311	
MS-Diffusion	0.336	0.354	0.335	0.339	0.342	0.322	0.339	0.343	0.349	0.332	0.333	0.338	0.337	0.344	0.345	0.324	0.321	0.335	0.344	0.347	0.349	0.353	0.352	0.331	0.344	0.327	0.319	0.317	0.324	
OminiControl	0.328	0.345	0.333	0.337	0.335	0.317	0.337	0.336	0.351	0.331	0.327	0.333	0.328	0.337	0.339	0.331	0.331	0.327	0.336	0.342	0.341	0.347	0.349	0.328	0.344	0.321	0.317	0.321	0.322	
SSR-Encoder	0.290	0.315	0.289	0.301	0.295	0.291	0.291	0.311	0.310	0.294	0.297	0.291	0.306	0.304	0.311	0.300	0.298	0.304	0.307	0.299	0.309	0.302	0.281	0.309	0.293	0.267	0.249	0.266		
UNO	0.322	0.343	0.330	0.328	0.325	0.312	0.326	0.331	0.349	0.322	0.319	0.321	0.321	0.330	0.330	0.320	0.320	0.315	0.328	0.336	0.334	0.333	0.340	0.324	0.324	0.321	0.295	0.299	0.297	
Emu2	0.310	0.316	0.307	0.303	0.309	0.285	0.308	0.310	0.316	0.316	0.316	0.297	0.310	0.310	0.311	0.319	0.300	0.296	0.299	0.305	0.316	0.308	0.315	0.318	0.304	0.326	0.309	0.292	0.276	0.266
RealCustom++	0.327	0.339	0.332	0.336	0.333	0.322	0.331	0.337	0.352	0.326	0.326	0.336	0.323	0.338	0.335	0.322	0.323	0.326	0.340	0.338	0.339	0.344	0.346	0.328	0.343	0.325	0.312	0.315	0.309	
OminiGen	0.318	0.318	0.312	0.333	0.312	0.307	0.316	0.313	0.357	0.316	0.314	0.316	0.310	0.322	0.330	0.316	0.325	0.317	0.314	0.327	0.327	0.328	0.340	0.311	0.321	0.317	0.324	0.317	0.315	
Custom Diffusion	0.331	0.330	0.316	0.322	0.321	0.313	0.323	0.321	0.337	0.327	0.312	0.316	0.316	0.325	0.325	0.322	0.322	0.326	0.325	0.329	0.326	0.332	0.336	0.317	0.330	0.310	0.312	0.314	0.315	
DreamBooth	0.315	0.330	0.315	0.321	0.320	0.314	0.319	0.325	0.299	0.316	0.312	0.306	0.318	0.320	0.324	0.313	0.320	0.321	0.320	0.327	0.325	0.330	0.323	0.303	0.326	0.301	0.319	0.319	0.316	
Textual Inversion	0.310	0.307	0.294	0.308	0.304	0.275	0.295	0.306	0.323	0.312	0.296	0.307	0.300	0.292	0.309	0.292	0.302	0.294	0.299	0.305	0.307	0.320	0.301	0.316	0.283	0.276	0.290	0.286	0.286	
$\lambda$ -Eclipse	0.280	0.282	0.282	0.296	0.296	0.298	0.278	0.297	0.312	0.278	0.290	0.293	0.313	0.310	0.315	0.284	0.289	0.289	0.292	0.301	0.290	0.297	0.301	0.294	0.287	0.272	0.274	0.248	0.248	
HiPer	0.313	0.335	0.316	0.316	0.316	0.304	0.316	0.316	0.308	0.312	0.311	0.307	0.311	0.307	0.319	0.319	0.313	0.317	0.301	0.326	0.331	0.333	0.324	0.324	0.313	0.320	0.305	0.305	0.	

Table 14: We evaluated the performance of various methods on DSH-Bench dataset, specifically analyzing their effectiveness across the third-level categories (under non-photorealistic).

Method	Subject Preservation																													
	VE	MI	PUF	FB	MS	BO	FUR	HA	AM	BU	DP	IN	ST	DN	PL	JE	BS	AR	CL	SE	SBA	TO	MA	RE	BI	FI	HF	FA	AC	
BLIP-Diffusion	0.227	0.170	0.175	0.208	0.146	0.153	0.210	0.180	0.133	0.185	0.139	0.175	0.108	0.208	0.183	0.171	0.181	0.125	0.188	0.375	0.221	0.177	0.219	0.247	0.192	0.183	0.180	0.221	0.194	
IP-Adapter	0.240	0.248	0.175	0.283	0.212	0.281	0.178	0.210	0.208	0.223	0.164	0.188	0.192	0.248	0.183	0.217	0.227	0.175	0.250	0.683	0.243	0.144	0.250	0.233	0.245	0.208	0.177	0.200	0.201	
MS-Diffusion	0.362	0.410	0.387	0.375	0.396	0.361	0.363	0.345	0.342	0.292	0.272	0.292	0.233	0.438	0.250	0.296	0.310	0.367	0.488	0.742	0.426	0.235	0.325	0.342	0.363	0.300	0.302	0.233	0.300	
OminControl	0.290	0.370	0.367	0.333	0.377	0.275	0.181	0.252	0.292	0.265	0.269	0.213	0.237	0.342	0.075	0.254	0.362	0.217	0.300	0.633	0.342	0.217	0.273	0.261	0.190	0.283	0.158	0.092	0.202	
SSR-Encoder	0.187	0.168	0.221	0.260	0.183	0.219	0.179	0.198	0.208	0.201	0.103	0.079	0.146	0.173	0.179	0.133	0.154	0.175	0.170	0.667	0.260	0.131	0.219	0.156	0.177	0.208	0.145	0.208	0.176	
UNO	0.440	0.448	0.537	0.365	0.467	0.428	0.415	0.399	0.433	0.377	0.267	0.308	0.342	0.471	0.317	0.462	0.481	0.458	0.530	0.675	0.438	0.335	0.350	0.431	0.358	0.375	0.315	0.208	0.338	
Emu2	0.350	0.310	0.337	0.315	0.412	0.303	0.303	0.257	0.192	0.267	0.186	0.122	0.379	0.448	0.100	0.333	0.333	0.158	0.305	0.433	0.350	0.187	0.294	0.314	0.280	0.375	0.251	0.250	0.288	
RealCustom++	0.435	0.425	0.763	0.431	0.508	0.472	0.436	0.380	0.525	0.300	0.331	0.404	0.208	0.469	0.337	0.646	0.398	0.433	0.317	0.450	0.493	0.342	0.419	0.406	0.392	0.367	0.265	0.179	0.271	
OminGen	0.235	0.130	0.108	0.246	0.348	0.231	0.232	0.180	0.292	0.152	0.094	0.200	0.158	0.175	0.067	0.158	0.256	0.100	0.243	0.375	0.169	0.083	0.263	0.242	0.253	0.183	0.178	0.142	0.231	
Custom Diffusion	0.071	0.072	0.067	0.048	0.073	0.094	0.017	0.052	0.033	0.154	0.019	0.037	0.050	0.077	0.050	0.029	0.038	0.060	0.047	0.150	0.042	0.002	0.042	0.017	0.055	0.092	0.225	0.000	0.056	
DreamBooth	0.171	0.245	0.179	0.215	0.181	0.211	0.111	0.170	0.233	0.208	0.089	0.096	0.163	0.231	0.050	0.196	0.154	0.050	0.172	0.308	0.112	0.040	0.120	0.089	0.118	0.233	0.065	0.017	0.087	
Textual Inversion	0.144	0.090	0.108	0.090	0.104	0.067	0.033	0.052	0.075	0.164	0.033	0.083	0.063	0.054	0.054	0.058	0.083	0.058	0.058	0.267	0.040	0.040	0.112	0.097	0.122	0.058	0.092	0.050	0.122	
A-Eclipse	0.279	0.237	0.225	0.317	0.244	0.261	0.194	0.218	0.250	0.226	0.214	0.192	0.092	0.235	0.221	0.192	0.219	0.450	0.227	0.517	0.250	0.196	0.270	0.294	0.232	0.267	0.215	0.196	0.234	
HiPer	0.162	0.167	0.165	0.192	0.142	0.142	0.057	0.114	0.108	0.162	0.053	0.079	0.075	0.156	0.054	0.067	0.153	0.200	0.135	0.333	0.126	0.021	0.109	0.031	0.107	0.183	0.061	0.046	0.108	
NEIT	0.181	0.200	0.237	0.146	0.194	0.164	0.189	0.150	0.275	0.181	0.142	0.154	0.117	0.210	0.146	0.125	0.148	0.317	0.093	0.192	0.135	0.177	0.210	0.231	0.182	0.150	0.301	0.285	0.227	0.254

Method	Prompt Following																												
	VE	MI	PUF	FB	MS	BO	FUR	HA	AM	BU	DP	IN	ST	DN	PL	JE	BS	AR	CL	SE	SBA	TO	MA	RE	BI	FI	HF	FA	AC
BLIP-Diffusion	0.290	0.287	0.291	0.293	0.291	0.279	0.275	0.284	0.281	0.280	0.292	0.277	0.278	0.296	0.259	0.282	0.287	0.259	0.286	0.285	0.292	0.277	0.283	0.291	0.274	0.289	0.256	0.222	0.272
IP-Adapter	0.317	0.336	0.316	0.321	0.335	0.304	0.307	0.319	0.326	0.304	0.313	0.320	0.321	0.332	0.267	0.320	0.331	0.293	0.320	0.319	0.320	0.312	0.321	0.317	0.311	0.309	0.314	0.308	0.327
MS-Diffusion	0.337	0.355	0.334	0.339	0.354	0.326	0.335	0.340	0.332	0.326	0.332	0.329	0.328	0.348	0.310	0.341	0.344	0.322	0.346	0.342	0.344	0.321	0.344	0.332	0.333	0.334	0.328	0.324	0.338
OminControl	0.330	0.348	0.341	0.335	0.340	0.321	0.334	0.333	0.341	0.330	0.324	0.334	0.330	0.344	0.321	0.329	0.336	0.309	0.339	0.334	0.334	0.326	0.338	0.328	0.342	0.333	0.326	0.313	0.337
SSR-Encoder	0.289	0.317	0.287	0.304	0.305	0.287	0.291	0.305	0.301	0.285	0.303	0.302	0.298	0.305	0.270	0.292	0.308	0.291	0.304	0.298	0.299	0.293	0.300	0.303	0.302	0.302	0.288	0.251	0.294
UNO	0.316	0.343	0.319	0.325	0.334	0.305	0.321	0.324	0.314	0.309	0.313	0.327	0.302	0.338	0.289	0.318	0.333	0.294	0.333	0.316	0.330	0.300	0.329	0.317	0.327	0.309	0.308	0.297	0.317
Emu2	0.309	0.337	0.303	0.311	0.326	0.303	0.293	0.305	0.313	0.298	0.307	0.317	0.301	0.318	0.291	0.309	0.304	0.276	0.308	0.323	0.312	0.273	0.297	0.312	0.304	0.297	0.292	0.243	0.279
RealCustom++	0.319	0.351	0.326	0.328	0.340	0.328	0.331	0.331	0.317	0.326	0.309	0.332	0.320	0.350	0.301	0.325	0.299	0.298	0.346	0.336	0.338	0.317	0.338	0.319	0.339	0.333	0.330	0.321	0.338
OminGen	0.322	0.321	0.313	0.328	0.341	0.318	0.308	0.305	0.332	0.317	0.301	0.322	0.309	0.321	0.299	0.298	0.338	0.305	0.326	0.336	0.314	0.323	0.329	0.327	0.309	0.333	0.326	0.315	0.333
Custom Diffusion	0.315	0.336	0.330	0.326	0.326	0.313	0.324	0.323	0.322	0.323	0.308	0.302	0.311	0.326	0.314	0.324	0.323	0.318	0.321	0.329	0.321	0.317	0.331	0.320	0.325	0.313	0.316	0.312	0.326
DreamBooth	0.322	0.337	0.329	0.331	0.330	0.317	0.322	0.317	0.327	0.322	0.310	0.317	0.309	0.322	0.314	0.319	0.316	0.295	0.325	0.334	0.325	0.325	0.334	0.328	0.323	0.303	0.324	0.315	0.324
Textual Inversion	0.290	0.348	0.341	0.335	0.340	0.321	0.334	0.333	0.341	0.330	0.324	0.334	0.330	0.344	0.321	0.329	0.336	0.309	0.339	0.334	0.334	0.326	0.338	0.328	0.342	0.333	0.326	0.313	0.337
A-Eclipse	0.280	0.313	0.256	0.300	0.293	0.290	0.295	0.292	0.301	0.278	0.311	0.304	0.302	0.310	0.273	0.293	0.306	0.241	0.311	0.335	0.304	0.288	0.305	0.297	0.306	0.296	0.286	0.283	0.279
HiPer	0.321	0.339	0.317	0.334	0.325	0.317	0.324	0.320	0.327	0.324	0.311	0.317	0.320	0.325	0.318	0.324	0.319	0.220	0.323	0.336	0.322	0.333	0.336	0.330	0.320	0.310	0.319	0.313	0.319
NEIT	0.297	0.327	0.305	0.306	0.317	0.292	0.274	0.291	0.302	0.295	0.272	0.300	0.301	0.305	0.257	0.286	0.300	0.256	0.291	0.325	0.288	0.270	0.299	0.290	0.295	0.301	0.285	0.227	0.254

Method	Image Quality																												
	VE	MI	PUF	FB	MS	BO	FUR	HA	AM	BU	DP	IN	ST	DN	PL	JE	BS	AR	CL	SE	SBA	TO	MA	RE	BI	FI	HF	FA	AC
BLIP-Diffusion	0.244	0.208	0.193	0.221	0.184	0.185	0.186	0.202	0.232	0.243	0.190	0.234	0.220	0.221	0.189	0.224	0.220	0.204	0.212	0.184	0.225	0.189	0.246	0.272	0.229	0.256	0.224	0.216	0.235
IP-Adapter	0.285	0.255	0.239	0.263	0.230	0.231	0.227	0.245	0.309	0.279	0.245	0.271	0.259	0.258	0.229	0.278	0.264	0.242	0.240	0.229	0.251	0.227	0.284	0.310	0.263	0.311	0.283	0.293	0.305
MS-Diffusion	0.304	0.287	0.268	0.296	0.269	0.273	0.269	0.276	0.308	0.264	0.267	0.286	0.288	0.294	0.267	0.299	0.298	0.267	0.284	0.283	0.289	0.255	0.307	0.321	0.296	0.332	0.293	0.307	0.311
OminControl	0.295	0.290	0.265	0.287	0.269	0.274	0.278	0.278	0.314	0.298	0.277	0.281	0.284	0.295	0.288	0.298	0.285	0.265	0.282	0.281	0.287	0.270	0.312	0.309	0.307	0.325	0.297	0.301	0.309
SSR-Encoder	0.258	0.237	0.216	0.241	0.204	0.206	0.213	0.239	0.288	0.260	0.214	0.242	0.232	0.234	0.217	0.245	0.237	0.227	0.229	0.210	0.241	0.212	0.271	0.287	0.266	0.293	0.260	0.248	0.268
UNO	0.286	0.277	0.249	0.288	0.257	0.249	0.252	0.262	0.312	0.275	0.243	0.285	0.267	0.290	0.223	0.282	0.274	0.226	0.265	0.228	0.276	0.240	0.302	0.315	0.290	0.311	0.273	0.278	0.294
Emu2	0.278	0.261	0.239	0.259	0.236	0.239	0.227	0.235	0.294	0.269	0.242	0.278	0.246	0.265	0.254	0.252	0.252	0.247	0.244	0.255	0.254	0.223	0.269	0.296	0.275	0.302	0.268	0.256	0.264
RealCustom++	0.304	0.300	0.283	0.291	0.269	0.280	0.277	0.274	0.316	0.314	0.258	0.304	0.290	0.304	0.250	0.296	0.282	0.277	0.306	0.284	0.300	0.259	0.311	0.331	0.313	0.321	0.305	0.319	0.330
OminGen	0.280	0.251	0.236	0.262																									

1296

1297

1298

1299

Table 15: We evaluated the performance of various methods on DSH-Bench dataset, specifically analyzing their effectiveness across **prompts with different scenarios**. Subject preservation, prompt following, and image quality are evaluated using SICS, CLIP-T, and HPSv2, respectively.

1300

1301

1302

1303

1304

1305

1306

1307

1308

1309

1310

1311

1312

1313

1314

1315

1316

1317

1318

1319

1320

1321

1322

1323

1324

1325

1326

1327

1328

1329

1330

1331

1332

1333

1334

1335

1336

1337

1338

1339

1340

1341

1342

1343

1344

1345

1346

1347

1348

1349

Subject Preservation						
Method	Background Change	Variation in Subject Viewpoint or Size	Interaction with Other Entities	Attribute Change	Style Change	Imagination
BLIP-Diffusion	0.204	0.207	0.201	0.182	0.189	0.195
IP-Adapter	0.233	0.230	0.224	0.177	0.203	0.209
MS-Diffusion	0.361	0.359	0.337	0.266	0.294	0.308
OminiControl	0.300	0.263	0.212	0.176	0.252	0.211
SSR-Encoder	0.206	0.201	0.200	0.166	0.171	0.188
UNO	<b>0.433</b>	<b>0.414</b>	<b>0.379</b>	<b>0.359</b>	<b>0.418</b>	<b>0.349</b>
Emu2	0.393	0.316	0.315	0.326	0.224	0.239
RealCustom++	0.386	0.384	0.353	0.297	0.314	0.310
OmniGen	0.238	0.167	0.159	0.125	0.155	0.133
Custom Diffusion	0.073	0.060	0.053	0.047	0.037	0.047
DreamBooth	0.180	0.157	0.138	0.139	0.128	0.144
Textual Inversion	0.121	0.102	0.104	0.109	0.074	0.098
$\lambda$ -Eclipse	0.262	0.257	0.249	0.246	0.244	0.230
HiPer	0.148	0.130	0.127	0.116	0.106	0.125
NeTI	0.211	0.182	0.185	0.198	0.173	0.182
<i>Aver.</i>	0.250	0.229	0.216	0.195	0.199	0.198
Prompt Following						
Method	Background Change	Variation in Subject Viewpoint or Size	Interaction with Other Entities	Attribute Change	Style Change	Imagination
BLIP-Diffusion	0.297	0.275	0.264	0.285	0.272	0.271
IP-Adapter	0.326	0.319	0.319	0.312	0.306	0.310
MS-Diffusion	0.342	<b>0.339</b>	<b>0.341</b>	0.324	<b>0.338</b>	0.341
OminiControl	0.338	0.334	0.337	<b>0.329</b>	0.326	<b>0.342</b>
SSR-Encoder	0.310	0.296	0.288	0.299	0.288	0.291
UNO	0.334	0.328	0.333	0.305	0.302	0.335
Emu2	0.308	0.305	0.297	0.295	0.313	0.308
RealCustom++	<b>0.343</b>	0.333	0.329	0.319	0.328	0.338
OmniGen	0.327	0.325	0.328	0.311	0.315	0.327
Custom Diffusion	0.326	0.317	0.320	0.328	0.325	0.323
DreamBooth	0.326	0.318	0.319	0.319	0.323	0.320
Textual Inversion	0.303	0.299	0.300	0.296	0.298	0.296
$\lambda$ -Eclipse	0.302	0.292	0.289	0.285	0.286	0.299
HiPer	0.325	0.323	0.315	0.318	0.318	0.311
NeTI	0.309	0.305	0.301	0.296	0.295	0.297
<i>Aver.</i>	0.321	0.314	0.312	0.308	0.309	0.314
Image Quality						
Method	Background Change	Variation in Subject Viewpoint or Size	Interaction with Other Entities	Attribute Change	Style Change	Imagination
BLIP-Diffusion	0.234	0.220	0.199	0.235	0.239	0.214
IP-Adapter	0.269	0.263	0.258	0.272	0.276	0.259
MS-Diffusion	0.291	0.292	0.292	0.287	0.300	0.301
OminiControl	0.285	0.283	0.290	<b>0.293</b>	0.294	0.296
SSR-Encoder	0.256	0.246	0.231	0.256	0.256	0.238
UNO	0.282	0.281	0.283	0.268	0.275	0.276
Emu2	0.262	0.260	0.249	0.252	0.268	0.270
RealCustom++	<b>0.300</b>	<b>0.298</b>	<b>0.295</b>	0.284	<b>0.301</b>	<b>0.307</b>
OmniGen	0.263	0.259	0.271	0.257	0.260	0.282
Custom Diffusion	0.245	0.236	0.237	0.248	0.231	0.240
DreamBooth	0.252	0.242	0.239	0.250	0.246	0.243
Textual Inversion	0.228	0.222	0.222	0.230	0.225	0.221
$\lambda$ -Eclipse	0.247	0.242	0.233	0.241	0.249	0.242
HiPer	0.250	0.247	0.241	0.255	0.247	0.240
NeTI	0.239	0.231	0.230	0.240	0.236	0.230
<i>Aver.</i>	0.260	0.255	0.251	0.258	0.260	0.257

- **Applicability:** Highly suitable for our scenario, especially when metrics have different scales or are not linearly related.

### 3. Kendall Rank Correlation Coefficient

1350  
1351  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
1360  
1361  
1362  
1363  
1364  
1365  
1366  
1367  
1368  
1369  
1370  
1371  
1372  
1373  
1374  
1375  
1376  
1377  
1378  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387  
1388  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1400  
1401  
1402  
1403

Table 16: We evaluated the performance of various methods on DSH-Bench dataset, specifically analyzing their effectiveness across **images with different difficulty levels**. Subject preservation, prompt following, and image quality are evaluated using SICCS, CLIP-T, and HPSv2, respectively.

Method	Subject Preservation			Prompt Following			Image Quality		
	Easy	Medium	Hard	Easy	Medium	Hard	Easy	Medium	Hard
BLIP-Diffusion	0.221	0.209	0.190	0.284	0.278	0.273	0.198	0.227	0.232
IP-Adapter	0.266	0.233	0.206	0.316	0.315	0.316	0.236	0.270	0.278
MS-Diffusion	0.410	0.362	0.312	<b>0.340</b>	<b>0.339</b>	<b>0.335</b>	0.278	0.297	0.299
OminiControl	0.294	0.256	0.242	0.337	0.336	0.331	0.278	0.292	0.294
SSR-Encoder	0.234	0.212	0.174	0.299	0.295	0.294	0.220	0.251	0.257
UNO	<b>0.469</b>	<b>0.405</b>	<b>0.383</b>	0.326	0.325	0.319	0.261	0.281	0.283
Emu2	0.349	0.346	0.332	0.308	0.306	0.301	0.239	0.263	0.268
RealCustom++	0.448	0.379	0.331	0.334	0.333	0.329	<b>0.281</b>	<b>0.300</b>	<b>0.303</b>
OmniGen	0.224	0.188	0.170	0.321	0.324	0.321	0.249	0.267	0.272
Custom Diffusion	0.067	0.061	0.060	0.323	0.324	0.322	0.234	0.241	0.241
DreamBooth	0.184	0.163	0.139	0.323	0.322	0.319	0.232	0.248	0.249
Textual Inversion	0.092	0.112	0.115	0.295	0.300	0.299	0.206	0.226	0.233
$\lambda$ -Eclipse	0.286	0.260	0.235	0.302	0.293	0.286	0.228	0.244	0.248
HiPer	0.139	0.145	0.122	0.323	0.319	0.315	0.230	0.251	0.251
NeTI	0.203	0.189	0.190	0.303	0.302	0.298	0.214	0.237	0.242
<i>Aver.</i>	0.259	0.235	0.213	0.316	0.314	0.311	0.239	0.260	0.263

Table 17: **Experiment hyperparameters on DSH-Bench**. LR: learning rate, Steps: training steps, GS: guidance scale

Method	T2I Model	Batch Size	LR	Train Steps	GS	Infer Steps	Additional parameter
BLIP-Diffusion	SD v1.5	N/A	N/A	N/A	7.5	25	N/A
IP-Adapter	SDXL	N/A	N/A	N/A	7.5	30	ip_adapter_scale: 0.5
MS-Diffusion	SDXL	N/A	N/A	N/A	7.5	30	scale: 0.6
OminiControl	FLUX.1-schnell	N/A	N/A	N/A	3.5	10	condition_scale: 1
SSR-Encoder	SD v1.5	N/A	N/A	N/A	7.5	30	$\lambda$ : 0.5
UNO	FLUX.1-dev	N/A	N/A	N/A	4	25	N/A
Emu2	SDXL	N/A	N/A	N/A	3.0	50	N/A
RealCustom++	SDXL	N/A	N/A	N/A	7.5	25	N/A
OmniGen	SDXL	N/A	N/A	N/A	2.5	50	img_guidance_scale: 1.8
$\lambda$ -Eclipse	SDXL	N/A	N/A	N/A	7.5	50	N/A
Textual Inversion	SD v1.5	4	5e-4	3000	7.5	50	N/A
DreamBooth	SD v1.5	1	2.5e-6	250	7.5	50	N/A
Custom Diffusion	SD v1.4	2	1e-5	250	6.0	100	N/A
HiPer	SD v1.4	1	5e-3	1500	7.5	50	N/A
NeTI	SD v1.4	2	1e-3	250	7.5	50	N/A

- **Advantages:** Also measures monotonic relationships; robust to outliers; suitable for rank/ordinal data.
- **Disadvantages:** More computationally intensive than Spearman; only captures monotonic relationships.
- **Applicability:** Also highly suitable, especially for smaller datasets or when we want a more robust rank-based measure.

#### 4. Krippendorff’s Alpha

- **Advantages:** Handles multiple raters and various data types (nominal, ordinal, interval, ratio); can handle missing data.
- **Disadvantages:** Mainly used for inter-rater reliability, not for correlation; does not indicate the direction of association; computationally complex.
- **Applicability:** Not suitable for our scenario, as it is designed to measure agreement among multiple raters.

Consequently, we choose *Kendall’s  $\tau$  value* and *Spearman correlation coefficient value*.

### E.3 DETAILS OF PROMPT GENERATION

We present below a comparative analysis of prompts generated by three state-of-the-art vision-language models: GPT-4o, Gemini 2.5 Pro, and Claude. The specific instruction used to generate the

---

1404 prompts can be found in Figure 5. A qualitative evaluation of these prompts reveals no significant  
1405 or discernible differences in their content or structure. Given the comparable performance across  
1406 the models, we selected GPT-4o for all prompt generation tasks in this study. Moreover, the results  
1407 generated by GPT-4o include slightly more complex and diverse descriptions of the environment,  
1408 which makes the model’s generation and application more intricate and challenging. This choice also  
1409 aligns with the methodology established in the Dreambench++, ensuring methodological consistency  
1410 with prior work. The example is as follows:

#### 1411 **Background Change**

- 1412
- 1413 1. **Gemini 2.5 Pro:** A single african drum resting on the red earth of the African savanna at  
1414 sunset.
- 1415 2. **Claude:** An african drum standing on a sandy beach with waves in the background
- 1416 3. **GPT-4o:** A single african drum placed on a sandy beach with gentle waves in the background  
1417 under a clear sky.
- 1418 4. **Gemini 2.5 Pro:** An african drum placed on a clean, white studio background.
- 1419 5. **Claude:** An african drum placed in a lush green rainforest clearing.
- 1420 6. **GPT-4o:** An african drum positioned on a wooden table inside a cozy rustic cabin with  
1421 warm lighting and wooden walls.
- 1422

#### 1423 **Variation in Subject Viewpoint or Size**

- 1424
- 1425 1. **Gemini 2.5 Pro:** A low-angle shot of an african drum sitting on a wooden stage, illuminated  
1426 by a single, warm spotlight from above, with dust motes dancing in the light.
- 1427 2. **Claude:** A low-angle shot of an african drum on a stage, with spotlights creating dramatic  
1428 shadows.
- 1429 3. **GPT-4o:** An african drum viewed from a high-angle perspective, placed on a grassy hilltop  
1430 overlooking a distant mountain range with scattered wildflowers around it.
- 1431 4. **Gemini 2.5 Pro:** Bird’s-eye view of an african drum next to a crackling campfire at dusk,  
1432 with the flickering firelight casting long, dancing shadows on the ground.
- 1433 5. **Claude:** A bird’s-eye view of an african drum in the center of a circular stone pattern in a  
1434 park.
- 1435 6. **GPT-4o:** A close-up shot of an african drum resting on a cobblestone pathway in a small  
1436 village, with faint rays of sunlight filtering through nearby trees and a soft shadow cast  
1437 beneath it.
- 1438

#### 1439 **Interaction with Other Entities**

- 1440
- 1441 1. **Gemini 2.5 Pro:** A close-up shot of a musician’s hands, blurred with motion, striking the  
1442 surface of an african drum, with beads of sweat flying off the drumhead.
- 1443 2. **Claude:** An african drum being played by multiple hands in a crowded marketplace, with  
1444 motion blur on the hands.
- 1445 3. **GPT-4o:** An african drum being played by a young boy surrounded by a group of people  
1446 dancing in a circle under the moonlight at an outdoor festival, motion blurs capturing the  
1447 dynamic movements of the dancers.
- 1448 4. **Gemini 2.5 Pro:** A group of children in a circle, their joyful faces illuminated by the sun,  
1449 all reaching out to play a large african drum together.
- 1450 5. **Claude:** An african drum partially submerged in a clear lake, with its reflection visible on  
1451 the water’s surface
- 1452 6. **GPT-4o:** An african drum partially obscured by tall reeds in a dense wetland, with a heron  
1453 standing nearby and its reflection visible in the water below.
- 1454

#### 1455 **Attribute Change**

- 1456
- 1457 1. **Gemini 2.5 Pro:** An african drum carved entirely from shimmering, translucent crystal,  
with glowing runes etched onto its surface.

- 
- 1458 2. **Claude:** An african drum made entirely of clear, shimmering crystal.
  - 1459
  - 1460 3. **GPT-4o:** An african drum painted entirely in vibrant rainbow colors and decorated with
  - 1461 intricate tribal patterns.
  - 1462 4. **Gemini 2.5 Pro:** A futuristic, metallic african drum with a holographic drumhead that
  - 1463 displays vibrant sound visualizations with every beat.
  - 1464 5. **Claude:** An african drum with a top made of stained glass, creating colorful patterns when
  - 1465 light shines through.
  - 1466 6. **GPT-4o:** An african drum crafted from polished metal instead of wood, with a reflective
  - 1467 surface that mirrors its surroundings.
  - 1468

### 1469 **Style Change**

- 1470
- 1471 1. **Gemini 2.5 Pro:** A cubist painting of an african drum, deconstructed into geometric shapes
- 1472 and shown from multiple perspectives simultaneously.
- 1473 2. **Claude:** An african drum rendered in the style of Vincent van Gogh's 'Starry Night', with
- 1474 swirling brushstrokes and vibrant colors.
- 1475 3. **GPT-4o:** An african drum depicted in the Cubist art style, with fragmented geometric shapes
- 1476 and abstract distortions.
- 1477 4. **Gemini 2.5 Pro:** An african drum depicted in the style of Japanese Ukiyo-e woodblock
- 1478 print, with flat colors, bold outlines, and a serene composition.
- 1479 5. **Claude:** A cubist interpretation of an african drum, broken down into geometric shapes and
- 1480 multiple perspectives.
- 1481 6. **GPT-4o:** An african drum illustrated in watercolor art, blending soft pastel tones with fluid
- 1482 brushstrokes, surrounded by splashes of color in the background.
- 1483
- 1484

### 1485 **Imagination**

- 1486
- 1487 1. **Gemini 2.5 Pro:** A fleet of miniature african drums floating through a cosmic nebula,
- 1488 propelled by rhythmic sound waves that ripple through the stardust.
- 1489 2. **Claude:** An african drum as a spaceship, with tiny alien creatures using it to explore the
- 1490 galaxy.
- 1491 3. **GPT-4o:** An african drum floating in mid-air, surrounded by glowing orbs of light that pulse
- 1492 rhythmically as if responding to the drum's silent beat.
- 1493 4. **Gemini 2.5 Pro:** In an enchanted forest, an ancient african drum is covered in moss and
- 1494 glowing mushrooms; when played, it causes the surrounding trees to grow and bloom
- 1495 instantly.
- 1496 5. **Claude:** An african drum transformed into a living creature, with eyes and limbs, dancing
- 1497 in a magical forest.
- 1498 6. **GPT-4o:** An african drum transformed into a magical portal, with swirling galaxies and
- 1499 stars emerging from its open top.
- 1500

## 1501 E.4 DETAILS OF HUMAN ANNOTATION

### 1502 E.4.1 ANNOTATION VERIFICATION GUIDELINES FOR SUBJECT DIFFICULTY LEVEL

#### 1503 CLASSIFICATION

#### 1504 **Task Objective**

- 1505
- 1506
- 1507
- 1508
  - Verify the difficulty label assigned by the model (Easy / Medium / Hard) for the “subject
  - 1509 detail preservation” when an image is used as a reference image.
  - Subject: the primary, most prominent object or semantic entity in the image
- 1510
- 1511

#### **Label Definitions (must align strictly)**

- 
- 1512
- **Easy:** Low surface complexity, homogeneous texture, near-uniform color/material; virtually no high-frequency details. Examples: smooth, solid-colored ceramic mug; smooth sphere; plain object without text/markings.
  - **Medium:** Contains discernible high-frequency features while maintaining globally simple, coherent structure; local details are present but not overwhelmingly dense. Examples: cylindrical container with readable text/logo; simple shapes with a few clear markings/scales/brand labels.
  - **Hard:** Non-uniform texture distribution with multi-scale geometric/details; dense, fine-grained features across regions that materially affect the subject's appearance. Examples: book cover with fine calligraphy and intricate patterns; woven/engraved/ornamented materials; many lines of small text with layout hierarchy.

### 1523 **Step-by-Step Decision Process**

#### 1524 **1. Identify the subject**

- Select the single most salient object (by size, focus sharpness, centrality, semantic importance).
- Multiple visible parts forming one entity (e.g., front and back of one book) count as one subject.

#### 1531 **2. Assess scale and detail density**

- High-frequency details: small text, granular textures, dense lines/patterns, fine edges, repetitive microstructures.
- Multi-scale: presence of large contours plus mid/small-scale text/patterns that meaningfully define appearance.

#### 1536 **3. Evaluate texture uniformity and surface complexity**

- Uniform materials (solid/near-solid color, smooth, gradual highlights) → tends to *Easy*.
- Simple shape with a few clear elements (e.g., a single line of text, logo) → tends to *Medium*.
- Complex surface patterns with non-uniform textures across regions, multiple areas with distinct details → often *Hard*.

#### 1543 **4. Evaluate readable elements**

- More numerous/smaller/denser text, more font variations, and deeper layout hierarchy increase difficulty.
- A single large word or single big logo alone does not imply *Hard*; typically *Medium*.

#### 1547 **5. Shape and geometry**

- Basic primitives (sphere/cylinder/cube) without details → *Easy*.
- Basic primitives with a few marks/engraved lines → *Medium*.
- Irregular shapes, folds/pleats, filigree, cutouts, elaborate ornaments → often *Hard*, depending on detail density.

#### 1552 **6. Lighting and material cues (auxiliary only)**

- Specular highlights or reflections alone do not equal high-frequency detail. If the surface lacks true microstructure/patterns, it can still be *Easy*.
- Genuine material textures (wood grain, fabric weave, leather grain, brushed metal) may raise difficulty to *Medium/Hard* based on density and multi-scale complexity.

#### 1558 **7. Context and occlusion (consider only if they affect visible details)**

- Complex backgrounds do not increase subject difficulty unless inseparable from or part of the subject surface.

### 1562 **Decision Boundaries and Edge Cases**

- *Easy vs. Medium:* Are there local details that must be faithfully reproduced to recognize the subject (e.g., small text, scales, granular texture)? If present and more than minimal, lean *Medium*.

- 
- 1566 • *Medium vs. Hard*: Are there multi-scale, regionally distributed complex details (e.g., title +
  - 1567 subtitle + fine print over a textured background; or patterns + filigree + material grain) that
  - 1568 jointly define appearance? If yes, lean *Hard*.
  - 1569 • **Text cases**:
  - 1570 – Single large brand word: *Medium*.
  - 1571 – Multiple lines of small text, varied font sizes, complex hierarchy: *Hard*.
  - 1572
  - 1573 • **Texture cases**:
  - 1574 – Uniform matte/gloss: *Easy*.
  - 1575 – Discernible but sparse texture (e.g., coarse weave with large spacing): *Medium*.
  - 1576 – Fine, dense, non-uniform texture (e.g., fine weave, leather grain, wood grain overprinted
  - 1577 with graphics): *Hard*.
  - 1578
  - 1579 • **Shape cases**:
  - 1580 – Smooth ceramic mug (no pattern): *Easy*.
  - 1581 – Measuring cup with a few scales/digits: *Medium*.
  - 1582 – Book cover with fine calligraphy plus intricate patterns/material grain: *Hard*.
  - 1583

1584 **Operational Protocol and Output Format** For each image, output:

- 1585 1. **Difficulty verification**: choose *Easy / Medium / Hard / Uncertain*.
- 1586 2. **Consistency with model label**: *Consistent / Inconsistent*.
- 1587

### 1588 **Quality Control and Consistency**

- 1589
- 1590 • Reduce borderline oscillation. If on the boundary, use majority voting over three axes:
- 1591 high-frequency detail density, texture uniformity, multi-scale complexity.
- 1592 • If image quality is too low (blur/low resolution) to assess detail density, mark *Uncertain* and
- 1593 state the reason.
- 1594 • Do not score based on aesthetics or unrelated factors (exposure, colorfulness), unless they
- 1595 impair detail discernibility.
- 1596

### 1597 E.4.2 ANNOTATION VERIFICATION GUIDELINES FOR PROMPT GENERATION

1598 This guidelines defines the validation criteria for human annotators to assess whether prompts

1599 generated by a large model, based on an image subject, comply with task standards and formatting

1600 requirements. Please verify each item strictly and record any non-compliance with brief explanations

1601 in the annotation tool.

1602

1603 **1. Overall Objective** Annotators should determine whether the model-generated prompts, given

1604 the clearly defined subject `subject_name` (provided by the task input or uniquely identifiable

1605 from the image), meet the following:

1606

- 1607 1. Content centers on the subject and matches the semantic requirements of the corresponding
- 1608 category (background change, variation in subject viewpoint or size, interaction with other
- 1609 entities, attribute changes, style changes, imagination).
- 1610 2. Quantity and structure meet the required counts.
- 1611 3. Language is clear, unambiguous, and executable as image-generation prompts.
- 1612 4. No prohibited or unsafe content (e.g., safety, ethics, copyright issues). Mark as non-
- 1613 compliant if present.
- 1614

### 1615 **2. Input and Output Specifications (Strict)**

- 1616
- 1617 • Input: an image (with a subject) and the model's JSON output.
- 1618 • The model output must be valid JSON with the following six top-level fields (names must
- 1619 match exactly):

- 
- 1620           – "background change": list with 2 items  
1621           – "variation in subject viewpoint or size": list with 2 items  
1622           – "interaction with other entities": list with 2 items  
1623           – "attribute changes": list with 2 items  
1624           – "style changes": list with 2 items  
1625           – "imagination": list with 2 items  
1626  
1627       • Each list element is a natural-language prompt (preferably a complete sentence) and must  
1628       explicitly include the subject name: `subject_name`.  
1629  
1630       • JSON must be machine-parseable: paired quotes, correct commas and brackets, and no  
1631       comments or trailing commas.

1632  
1633 **3. Category Definitions and Decision Criteria** Read each prompt and decide whether it satisfies  
1634 the semantic requirements of its assigned category. If not, label it “Non-compliant – Category Error.”

1635  
1636 (1) BACKGROUND CHANGE (2 ITEMS)

- 1637       • Requirement: Only change the background environment or setting; the subject’s default  
1638       attributes remain unchanged (color/material/shape/appearance unchanged). No interactions  
1639       or complex effects.  
1640       • Acceptable: Simple location/time/weather descriptions (e.g., in a park,” by the sea,” “at  
1641       sunset”).  
1642       • Prohibited: Large viewpoint changes (e.g., bird’s-eye view, low-angle close-up), strongly  
1643       stylized lighting techniques, obstruction/reflection/motion effects, interactions, or attribute  
1644       changes (e.g., color/material swaps).  
1645  
1646

1647 (2) VARIATION IN SUBJECT VIEWPOINT OR SIZE (2 ITEMS)

- 1648       • Requirement: Moderate adjustments to subject location, perspective, and lighting are  
1649       allowed; add non-interactive environmental elements (buildings, roads, trees, etc.).  
1650       • Must include more than lighting changes: at least one clear change in perspec-  
1651       tive/composition/scene elements.  
1652       • Prohibited: Interactions with other animals/people; obstruction, reflection, motion effects;  
1653       changes to inherent subject attributes.  
1654  
1655

1656 (3) INTERACTION WITH OTHER ENTITIES (2 ITEMS)

- 1657       • Requirement: Include one or more of the following:  
1658           – Interactions between the subject and other animals/people or multi-subject interactions.  
1659           – Visual effects such as obstruction, reflection, motion blur/trails.  
1660       • May combine perspective, lighting, and environmental changes.  
1661       • Prohibited: Changing the subject’s inherent attributes (color/shape/material/appearance); if  
1662       attributes change, it belongs in “attribute changes.”  
1663  
1664  
1665

1666 (4) ATTRIBUTE CHANGES (2 ITEMS)

- 1667       • Requirement: Change the subject’s own attributes (color, shape, material, components’  
1668       appearance, patterns, structural details).  
1669       • Changes should focus on the subject itself, not merely environment or lighting.  
1670       • Prohibited: Purely stylistic changes (e.g., in the style of Impressionism”) belong in style  
1671       changes”; scene/interaction-only changes belong in background change/variation in subject  
1672       viewpoint or size/interaction with other entities.  
1673

- 
- 1674 (5) STYLE CHANGES (2 ITEMS)  
1675  
1676 • Requirement: Change the overall artistic style (art movements, media, periods, processes,  
1677 rendering methods).  
1678 • Example dimensions: Impressionism, Cubism, Baroque, Minimalism, cyberpunk, water-  
1679 color, oil painting, pixel art, low-poly, woodcut/printmaking, cel shading, photorealism,  
1680 etc.  
1681 • Prohibited: Only changing inherent subject attributes (belongs in attribute changes”); only  
1682 changing perspective/interaction/obstruction (belongs in other categories).  
1683 • Copyright and sensitivity: Avoid using identifiable living/modern specific artist names or  
1684 protected signature styles. If present, mark Non-compliant – Copyright/Style.”  
1685

- 1686 (6) IMAGINATION (2 ITEMS)  
1687  
1688 • Requirement: Construct unrealistic or fantastical scenarios (science fiction, fantasy, surreal-  
1689 ism), allowing violations of physical laws or common sense.  
1690 • Must still center on the subject and include the subject name.  
1691 • Prohibited: Vague fantasies unrelated to the subject; avoid illegal, harmful, or clearly unsafe  
1692 content.  
1693  
1694

#### 1695 4. Common Compliance Requirements

##### 1696 SUBJECT CONSISTENCY 1697

- 1698  
1699 • Every prompt must contain the literal string `subject_name` (e.g., the subject is {sub-  
1700 ject\_name}” or subject\_name standing...”). If missing or misspelled, mark Non-compliant –  
1701 Subject Missing/Error.”  
1702 • Content must revolve around the subject, with a clear role and visibility. If the subject is  
1703 ignored or replaced, mark Non-compliant – Subject Irrelevant.”  
1704

##### 1705 SEMANTIC CLARITY AND NON-CONTRADICTION 1706

- 1707 • No contradictory descriptions (e.g., midnight with strong sunlight,” underwater and on a  
1708 desert” simultaneously).  
1709 • Do not include elements exceeding the category scope (e.g., complex interactions/effects  
1710 in background change/variation in subject viewpoint or size; attribute/style changes in  
1711 background change/variation in subject viewpoint or size/interaction with other entities).  
1712 • In attribute changes,” do not let style terms be the primary change; conversely, in style  
1713 changes,” avoid attribute edits as the main point.  
1714

##### 1715 DIVERSITY AND NON-REPETITION 1716

- 1717 • Prompts within the same field should have substantive differences (locat-  
1718 ion/viewpoint/elements/interaction types/attribute or style dimensions).  
1719 • Avoid mere synonym swaps or repeating the same composition.  
1720

##### 1721 SAFETY AND ETHICS 1722

- 1723 • No violent harm, hate, explicit adult content, illegal activity, or instructions for dangerous  
1724 behavior.  
1725 • For people, avoid identifiable personal or sensitive information.  
1726 • Avoid generating protected brand logos or named-artist signature styles. Prefer generic style  
1727 descriptors.

---

### 1728 E.4.3 ANNOTATION GUIDELINES FOR LABELING OF THE SICS TRAINING DATASETS

1729  
1730 For detailed annotation guidelines, please refer to the Appendix E.2.

### 1731 1732 E.4.4 ANNOTATOR RECRUITMENT, TRAINING, AND CONSENSUS PROTOCOL

1733 We provide the following elaboration on annotator requirements and the specifics of our annotation  
1734 pipeline:

- 1735  
1736 1. All annotators involved in this work possess extensive experience in domain-relevant anno-  
1737 tation. They have previously participated in similar tasks and have a thorough understanding  
1738 of subject-driven text-to-image (T2I) generation.
- 1739 2. All annotators hold a bachelor’s degree or higher and are between 25 and 30 years old,  
1740 ensuring the capability to accurately comprehend and implement the annotation guidelines.
- 1741 3. We provided human annotators with sufficient training to ensure full understanding of the  
1742 subject-driven T2I task and to promote unbiased, accurate evaluations. Detailed annotation  
1743 guidelines are provided (see Figure 5 and Figure 12), and we conducted meetings with all  
1744 annotators to comprehensively align and explain the rules.
- 1745 4. To ensure consensus on annotation standards, we conducted a calibration process. Annotators  
1746 first labeled a small pilot dataset; we then reviewed discrepant cases and provided targeted  
1747 feedback and additional training. This iterative process was repeated until all annotators  
1748 demonstrated high consistency in their understanding of the differentiation criteria for each  
1749 annotation task.
- 1750 5. During the formal annotation phase, each sample was independently labeled by five annota-  
1751 tors. To construct a high-confidence training dataset, we applied a consensus-based filtering  
1752 criterion, retaining only samples for which at least four annotators (i.e.,  $\geq 80\%$ ) assigned  
1753 identical labels.

## 1754 1755 F MORE GENERATION EXAMPLES

1756  
1757 Figure 13 shows the generation examples of different methods across different difficulty levels.  
1758 Figure 14 shows the generation examples of different methods across different prompt scenarios.  
1759 The blue block highlights encoder-based methods, and the green block highlights fine-tuning-based  
1760 methods.

## 1761 1762 1763 G LARGE LANGUAGE MODEL USAGE STATEMENT

1764  
1765 In this work, we leverage large language models to assist human annotators during data labeling  
1766 and to generate prompts. Detailed procedures are provided in the Appendix E.3, Appendix E.4,  
1767 Section 2.1.1 and Section 2.1.3. We also employed large language models to aid in manuscript  
1768 preparation.

## 1769 1770 H LIMITATIONS

1771  
1772 DSH-Bench addresses the limitations of current subject-driven T2I generation benchmarks by pro-  
1773 viding a comprehensive and diverse dataset with 459 subject images and 5,508 prompts, covering  
1774 categories such as person, mammal, clothing, and so on. However, the benchmark is constrained  
1775 to 459 subject images. Increasing the number of test samples could enhance the credibility and  
1776 complexity of the evaluation. Additionally, we did not conduct a cross-analysis between subject  
1777 difficulty and prompt scenario. Despite meticulous manual reviews, some unintentional annotation  
1778 errors may still be present.

1779  
1780  
1781

1782  
1783  
1784  
1785  
1786  
1787  
1788  
1789  
1790  
1791  
1792  
1793  
1794  
1795  
1796  
1797  
1798  
1799  
1800  
1801  
1802  
1803  
1804  
1805  
1806  
1807  
1808  
1809  
1810  
1811  
1812  
1813  
1814  
1815  
1816  
1817  
1818  
1819  
1820  
1821  
1822  
1823  
1824  
1825  
1826  
1827  
1828  
1829  
1830  
1831  
1832  
1833  
1834  
1835

	Easy	Medium	Hard
<b>Reference Image</b>			
<b>BLIP-Diffusion</b>			
<b>IP-Adapter</b>			
<b>MS-Diffusion</b>			
<b>OminiControl</b>			
<b>SSR-Encoder</b>			
<b>UNO</b>			
<b>Emu2</b>			
<b>RealCustom++</b>			
<b>OmniGen</b>			
<b>λ-Eclipse</b>			
<b>HiPer</b>			
<b>NeTI</b>			
<b>Custom Diffusion</b>			
<b>DreamBooth</b>			
<b>Textual Inversion</b>			
<b>Prompt</b>	A wooden stool placed in the middle of a lush green garden, surrounded by blooming flowers and soft grass.	A high-angle shot of a medicine bottle placed on the sand at the edge of a serene beach with gentle waves approaching and seashells scattered nearby.	A plain beige bowl placed on a wooden table with a scenic countryside view in the background.
		A book titled 'A BOOK FULL OF HOPE' lying on a sandy beach with soft waves visible in the background.	
		A woman sitting on a bench in a city park during autumn, surrounded by falling leaves and vibrant shades of orange and yellow. Perspective: medium shot, eye-level.	
		A penguin viewed from a bird's-eye perspective standing on an ice floe with scattered pieces of ice floating in a deep blue ocean under twilight lighting.	
		A yellow alarm clock photographed from a bird's-eye view, placed on a messy work desk filled with scattered papers, pens, and a coffee cup.	
		Capture the canned health products from a bird's-eye view in an urban garden, surrounded by vibrant foliage and colorful flowers. The lighting transitions through patches of sunlight and shadows created by trees.	
		A soccer ball resting on a sandy beach with waves gently rolling in the background, keeping the soccer ball's default attributes unchanged.	

Figure 13: Examples of images generated by all methods on different subjects difficulty level.

1836  
1837  
1838  
1839  
1840  
1841  
1842  
1843  
1844  
1845  
1846  
1847  
1848  
1849  
1850  
1851  
1852  
1853  
1854  
1855  
1856  
1857  
1858  
1859  
1860  
1861  
1862  
1863  
1864  
1865  
1866  
1867  
1868  
1869  
1870  
1871  
1872  
1873  
1874  
1875  
1876  
1877  
1878  
1879  
1880  
1881  
1882  
1883  
1884  
1885  
1886  
1887  
1888  
1889

	Background Change	Variation in Subject Viewpoint or Size	Interaction with Other Entities	Attribute Change	Style Change	Imagination
<b>Reference Image</b>						
<b>BLIP-Diffusion</b>						
<b>IP-Adapter</b>						
<b>MS-Diffusion</b>						
<b>OminiControl</b>						
<b>SSR-Encoder</b>						
<b>UNO</b>						
<b>Emu2</b>						
<b>RealCustom++</b>						
<b>OmniGen</b>						
<b><math>\lambda</math>-Eclipse</b>						
<b>HiPer</b>						
<b>NeTI</b>						
<b>Custom Diffusion</b>						
<b>DreamBooth</b>						
<b>Textual Inversion</b>						
<b>Prompt</b>	The green velvet sofa situated outdoors in a garden setting, surrounded by blooming flowers and lush greenery, maintaining the same sofa attributes.	A puppy lying under a tree in a park, captured from a higher angle shot, with colorful leaves around it and soft sunlight filtering through the branches.	A construction worker holding a hammer while working on a wooden structure, with a nail being driven into the wood, capturing the speed and dynamics of the hammer's movement.	Sneakers redesigned with a metallic reflective surface, giving them a futuristic and glowing look.	An elephant depicted in the style of Van Gogh's 'Starry Night', featuring swirling blue and yellow tones, with the elephant's body mirroring the swirling patterns of the sky.	The Temple of Heaven floating on a misty cloud island in the sky, with traditional architecture and a pagoda, with a waterfall cascading off the edge into the void below.

Figure 14: Examples of images generated by all methods on different prompt scenarios.

1890  
 1891  
 1892  
 1893  
 1894  
 1895  
 1896  
 1897  
 1898  
 1899  
 1900  
 1901  
 1902  
 1903  
 1904  
 1905  
 1906  
 1907  
 1908  
 1909  
 1910  
 1911  
 1912  
 1913  
 1914  
 1915  
 1916  
 1917  
 1918  
 1919  
 1920  
 1921  
 1922  
 1923  
 1924  
 1925  
 1926  
 1927  
 1928  
 1929  
 1930  
 1931  
 1932  
 1933  
 1934  
 1935  
 1936  
 1937  
 1938  
 1939  
 1940  
 1941  
 1942  
 1943

<p><b>Nano-Banana</b></p>	<p>The boat in motion, reflecting the sky and sunlight on its wet surface, with dolphins playfully swimming alongside it in the ocean.</p>	<p>Steam locomotive viewed from a low-angle shot as it exits a dark tunnel, surrounded by lush greenery and an overhanging cliff.</p>	<p>Place the edible oil bottle beside a chef prepping food, with oil drops mid-pour creating a motion effect, and an obstruction created by utensils partially in the foreground, giving depth.</p>
			
	<p><b>GPT-4o</b></p>	<p>A vintage camera placed beside a reflective puddle, showing a distorted mirror image against the backdrop of people walking along a rainy street.</p>	<p>A motorbike toy captured in a bird's-eye view resting on an urban sidewalk with faint chalk drawings and a warm evening streetlight casting gentle shadows.</p>
			

Figure 15: Examples of badcases generated by GPT-4o and Nano-Banana on hard examples