DSH-Bench: A Difficulty- and Scenario-Aware Benchmark with Hierarchical Subject Taxonomy for Subject-Driven Text-to-Image Generation

Anonymous Author(s)

Affiliation Address email

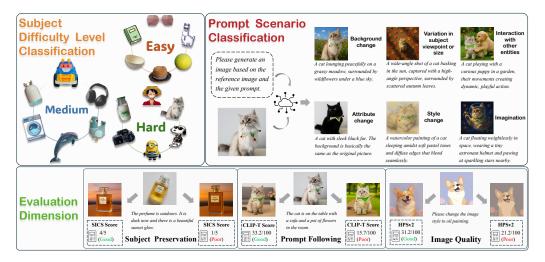


Figure 1: **Overview of DSH-Bench**. We curate a diverse dataset of subject images and categorize them into three difficulty levels—**easy**, **medium**, and **hard**—based on the complexity of preserving subject details. Leveraging GPT-4o's capabilities, we systematically generate contextually appropriate 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

5

6

8

9

10

11

12

13

14

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

28

29

30

31

32

34

35

36

37

38

42

43

44

45

46

47

48

49

50

51

52

53 54

55

56 57

58

59

60

61

62

63

64

65

66

67

Subject-driven text-to-image (T2I) generation aims to generate images conditioned on both textual 19 prompts and specific reference images. It has become feasible due to significant advancements in 20 large-scale T2I generative models [9, 13, 51, 47, 3, 5, 25, 10]. In subject-driven T2I generation, 21 aside from image quality considerations, two other fundamental criteria must be satisfied: Subject 22 Preservation and Prompt Following. Subject Preservation requires that the generated image accurately 23 maintain the details of the reference subject. Prompt Following demands that the generated image 24 consistently reflects the content in the prompt. For example, a user might request an image of "his dog 25 traveling around the world" [50]. In this scenario, the generated image must depict a dog identical to the reference image while illustrating the act of traveling as described in the prompt. 27

Significant progress has been made in subject-driven T2I generation in recent years [50, 14, 28, 58, 30, 70, 16, 62, 21, 45]. 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 [50, 28, 6, 59, 41] 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 2, 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 the aforementioned 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 [32], ImageNet [8], and category lists from Wikipedia [63] in the hierarchical taxonomy construction. As shown in Figure 3(a), the widely used DreamBench [50] 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× and 15×, respectively. Although DreamBench++ [41] 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 2 shows the model's performance varies significantly with different samples, highlighting the necessity for a classification of both subject image and prompt. Although DreamBench++ [41] categorizes prompts based on their perceived difficulty, the criteria underlying this classification are not clearly defined. Additionally, DreamBench++ [41] does not analyze the difficulty levels associated with different subjects. To address these limitations, we propose an innovative classification scheme. 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 [46] and DINO [4] with GPT-40 [37] for evaluation, resulting in improved alignment with

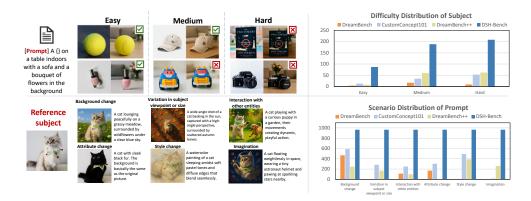


Figure 2: Qualitative comparison of generated images under different difficulty levels and scenarios.

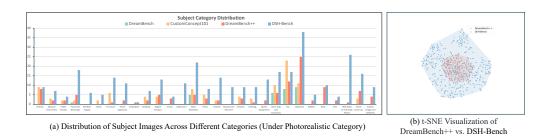


Figure 3: Distribution of subject images. (a) The distribution of images across different categories for DreamBench, CustomConcept101, DreamBench++, and DSH-Bench. (b) Images comparison between DSH-Bench and DreamBench++ using t-SNE

human evaluation. However, our benchmark reveals that per-model evaluation under this paradigm requires approximately 20,000 API calls to GPT-40, incurring prohibitive computational costs exceeding \$400 for each evaluation. To address the limitation, we introduce **Subject Identity Consistency Score** (SICS). 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 [2] on this dataset. Finally, we use Kendall's τ value to quantify the alignment between model outputs and human evaluation. Experimental results demonstrate that SICS achieves a statistically significant improvement, outperforming GPT-40 by 9.4% in human evaluation correlation metrics.

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 GPT-4o-based approaches. We are open-sourcing DSH-Bench, including all subject images, prompts, generated images, related code, and the SICS model.

2 Related Work

93

109

120

128

133

2.1 Subject-Driven Text-to-Image Generation

In recent years, subject-driven T2I generation has attracted significant research attention [50, 14, 94 28, 58, 30, 70, 16, 15, 62, 21, 45]. Within the context of diffusion models, optimization-based 95 model [14, 50, 28, 57, 34, 22, 18] enables subject-driven generation by introducing lightweight parameters and performs parameter-efficient fine-tuning for each subject. In contrast, the encoderbased methods [62, 70, 52, 35, 7, 31, 29, 49, 71, 20, 23, 67, 38, 64, 24, 19] leverage additional 98 image encoders and network layers to encode the reference image of the subject. ELITE [62] uses a learning-based encoder for subject customization, which consists of a global mapping network to 100 encode reference subjects into pseudo words and a local mapping network to maintain subject details. IP-Adapter [70] introduces cross-attention through an additional image encoder to incorporate control signals. Furthermore, SSR-Encoder [73] enhances identity preservation. This strategy facilitates subject-driven generation without necessitating further fine-tuning when introducing new concepts. 105 The Diffusion Transformers (DiT) [40] uses transformer as a denoising network to iteratively refine noisy image tokens, applied in T2I models widely [43, 48]. Based on these foundation models, 106 approaches like OminiControl [55] and UNO [64] explore the inherent image reference capabilities 107 of transformers, suggesting that DiT itself can serve as an image encoder for subject reference. 108

2.2 Subject-Driven T2I Generation Benchmark

Evaluation for subject-driven T2I generation involves a variety of metrics focusing on different aspects. For image quality, several notable studies [68, 27, 65, 1, 69, 60] have conducted Dream-Sim [12], CLIP-I [46], and DINO Score [4] are commonly adopted to measure perceptual similarity. 112 In terms of semantic consistency, the CLIP score [46] is frequently used. However, in subject-driven 113 image generation tasks, existing perceptual similarity metrics often diverge from human perception. 114 To address this limitation, researchers have proposed new metrics [41] that better align with human 115 judgments. DreamBench [50] is limited in the diversity of subjects and prompt scenarios. Dream-116 Bench++ [41] increases to 150 subject images. Moreover, current benchmarks can not provide a 117 systematic categorization of subjects and prompts, making it difficult to derive meaningful insights 118 from the evaluation results.

2.3 Subject Preservation Evaluation

Subject preservation evaluation plays a crucial role in the evaluation of subject-driven T2I generation.
Learning-based metrics [11, 72, 44] compute the distances between image features extracted by deep
neural networks. However, these approaches fall short in capturing the full range of nuances present
in human perception. To address this limitation, image embeddings from large vision models like
CLIP [46] and DINO [4] have been utilized. The image-retrieval score [33] has been used to assess
the visual similarity. To better align with human perceptual judgments, DreamSim [12] has been
introduced to assess image similarity with a focus on foreground objects and semantic content.

3 DSH-Bench

This section provides an overview of the primary components of DSH-Bench. Section 3.1 outlines the data construction process. In Section 3.2, we present a concise introduction to the definitions and evaluation methods for three evaluation dimensions. *A detailed explanation is available in the supplementary materials.*

3.1 Benchmark Dataset Construction

134 3.1.1 Subject Image Collection

Hierarchical Taxonomy Establishment As shown in Figure 4, we establish a hierarchical taxonomy. For the first- and second-level categories, we primarily refer to existing benchmarks from prior studies [50, 28, 41], resulting in two first-level categories and six second-level categories. For the third-level categories, we first reference COCO [32], ImageNet [8] and Wikipedia to compile a list of

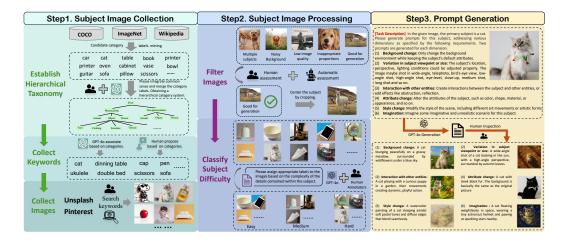


Figure 4: **Dataset construction process of DSH-Bench.** We construct a hierarchical taxonomy to obtain a comprehensive set of keywords. Then we collect web images using these keywords. After performing both manual review and automated filtering of the images, we classify the difficulty of subject images and use GPT-40 to generate prompts for each subject image.

candidate category labels, then utilize GPT-40 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++ [41], keywords collection relies on GPT-40 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-40 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—significantly surpassing DreamBench++'s 300. The specific keywords are provided in the Appendix B. Given a set of selected keywords, we retrieve images from Unsplash [56] and Pinterest [42]. Keywords without suitable images are discarded. We also add some excellent images from previous work. *Each image's copyright status has been verified for academic suitability*.

3.1.2 Subject Image Processing

Image Filtering To filter unsuitable images, human annotators remove images with multiple subjects and noisy backgrounds. We use aesthetic score [69] and SAM [26] 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 2, 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 instances 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 instances expose model limitations in maintaining structural fidelity and textural granularity under complex topological constraints. We utilize GPT-40 to classify the subject images according to the aforementioned criteria. Subsequently, all images are reviewed and corrected by human annotators to ensure accuracy and consistency.



Figure 5: Examples generated by methods listed in the leaderboard. Best viewed when zoomed in.

3.1.3 Prompt Generation

Although DreamBench++ categorizes prompts based on their perceived difficulty, it does not provide 173 empirical evidence to substantiate the criterion. To address this limitation, we organize the prompts 174 according to specific application scenarios, dividing them into six categories, including (1) Back-175 ground change (BC): scenarios involving changes in background elements. (2) Variation in subject 176 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 179 adherence to physical plausibility. (4) Attribute change (AC): scenarios involving modifications to 180 certain attributes of the subject, such as color or shape. (5) Style change (SC): scenarios involving 181 alterations in the artistic or visual style of the subject. (6) **Imagination (IM):** scenarios where the 182 target image depicts an imagined or fictional scene. We generate two prompts for each scenario. 183 The specific instructions employed for prompt generation are depicted in Figure 4. All prompts are 184 reviewed by two human annotators to ensure they are ethical and free from defects. 185

Finally, we obtain a total of **459** high-quality images and **5,508** prompts. Figure 2 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 3(b). The results clearly indicate that our benchmark achieves superior diversity.

3.2 Evaluation Dimension

190

191

192

193

194

196

197

198

199 200

201

202

203

204

205

206

207

Previous notable works [50, 14, 28, 58] evaluate the performance of subject-driven T2I models from two perspectives: Subject Preservation and Prompt Following. Mao et al. [36] also uses ImageReward [68] to evaluate image quality. Therefore, DSH-Bench evaluates from the three aforementioned dimensions.

Subject Preservation DreamBench++ [41] 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 [61] 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. 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 [46] is highly consistent with human annotations. Therefore, we also adopt CLIP-T score as the evaluation metric for prompt following.

Image Quality HPSv2 [65] utilizes professionally annotated data to more accurately reflect human aesthetic preferences for generated images. Previous studies [54] 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 in this work.

214 4 Experiment

4.1 Experiment Setup

Implementation Details We conduct experiments on two mainstream approaches: *i) Finetuning-based:* 1) Textual Inversion(TI) [15], 2) DreamBooth [50], 3) Custom Diffusion [28], 4) Hiper [17], 5) NeTI [1]. *ii) Encoder-based:* 1) BLIP-Diffusion [30], 2) IP-Adapter [70], 3) MS-Diffusion [59], 4) Emu2 [53], 5) OminiControl [55], 6) SSR-Encoder [73], 7) RealCustom++ [36], 8) OmniGen [66], 9) λ -Eclipse [39], 10) UNO [64]. Our experiments are conducted using the official implementations to guarantee reliability and fairness. More details can be found in Appendix E.

Human Annotation Five human annotators label the training datasets for SICS. To assess the alignment between various evaluation metrics and human evaluation, the same group of annotators is tasked with labeling the ground truth for images generated by each method on the DSH-Bench dataset. We provide human annotators with sufficient training to ensure they fully understand the subject-driven T2I generation task and can provide unbiased and discriminating scores.

Table 1: The human alignment degree among different evaluation metrics, measured by **Kendall's** τ **value** and **Spearman correlation coefficient value**. H: Human, G: GPT-40, D: DINO, Dv2: DINOv2, CB: CLIP-B, CL: CLIP-L, S: SICS.

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

4.2 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 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 (Bold font: the maximum value in a row. An underline: the second highest value in a row), it significantly outperforms GPT-40 [41] by 9.37% (KDV) and 5.31% (SCV). This performance gap strongly suggests SICS's enhanced capability in modeling human evaluation. Notably, GPT-40 demonstrates 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,

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

Method	Sub	ject Preserva	ation	Pr	ompt Follow	ing	I	mage Qualit	y
Method	DB	DB++	HB	DB	DB++	HB	DB	DB++	HB
BLIP-Diffusion	0.229	0.216	0.204	0.291	0.278	0.277	0.267	0.254	0.223
IP-Adapter	0.230	0.244	0.229	0.321	0.318	0.315	0.291	0.296	0.266
MS-Diffusion	0.316	0.346	0.352	0.332	0.339	0.338	0.311	0.314	0.294
OminiControl	0.279	0.268	0.258	0.325	0.337	0.334	0.312	0.308	0.290
SSR-Encoder	0.231	0.202	0.202	0.290	0.287	0.295	0.273	0.270	0.247
UNO	0.409	0.410	0.409	0.317	0.322	0.323	0.304	0.297	0.278
Emu2	0.360	0.343	0.341	0.291	0.309	0.304	0.272	0.278	0.260
RealCustom++	0.377	0.380	0.375	0.325	0.329	0.332	0.316	0.314	0.298

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	0.409	0.323	0.278	0.252
RealCustom++	SDXL	0.375	0.332	0.294	0.251
MS-Diffusion	SDXL	0.352	0.338	0.294	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
λ -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

prompts requiring attribute change constitute 22.7%, which is higher than the 16.7% observed in DSH-Bench. Figure 6(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:

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

SP, PF, and IQ represent the scores for Subject Preservation, Prompt Following, and Image Quality, respectively. λ, γ, μ 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 ensures that a model must perform well across all evaluation dimensions to achieve a high overall assessment. We rank all models based on S_h scores. Table 3 shows the leaderboard. UNO demonstrates relatively strong overall performance. We attribute this improvement to the novel architectural design of UNO and the minimal yet effective modifications implemented in DiT.

5 Analysis

In this section, we conduct a detailed analysis of the performance of all methods based on the hierarchical category system, the subject difficulty level classification, and the prompt scenario classification. The results are as follows:

A scientific and comprehensive subject image sampling method is necessary Figure 6(c) and Figure 6(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

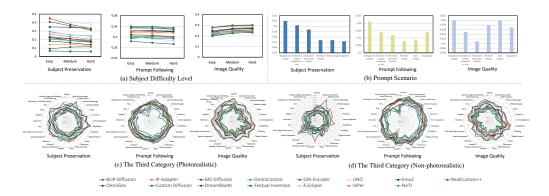


Figure 6: Comparison for DSH-Bench scores in different evaluation dimensions. The specific metric values are provided in the Appendix D.2. Best viewed when zoomed in.

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 6 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 6(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 6(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, as long as the generated image correctly represents the general category and overall shape, the evaluation score is unlikely to be substantially reduced, even if finer details are not perfectly captured. Given these findings, it is crucial to enhance models' ability to encode and reconstruct complex subject details more effectively in future research endeavors.

The subject-driven T2I capability for different prompt scenarios is not robust Figure 6(b) shows the average performance of all models across six prompt scenarios. The results show that: 1) In BC, VS, and IE scenarios, the model's performance consistently declines across all evaluation dimensions. This trend suggests that the difficulty of the scenarios increases progressively from BC to IE. Notably, the finding that the IE scenario is more challenging than the BC scenario aligns with intuitive expectations. 2) For subject preservation, the model's average performance across the AC, SC, and IM prompt scenarios remains relatively low. This could be because the generated subjects undergo partial modifications relative to the original subjects in these three scenarios. Given these findings, more emphasis should be placed on enhancing methods for IE prompt scenario. For instance, increasing the volume of training data tailored to these specific contexts.

6 Conclusion

This paper introduces a novel benchmark called DSH-Bench, designed specifically for subject-driven T2I generation. DSH-Bench presents unique challenges for subject-driven T2I generation models. Key features include: 1) a hierarchical category system in image collection to ensure both the diversity and comprehensiveness of subject images; 2) an innovative classification scheme for categorizing subject difficulty levels and prompt scenarios to obtain valuable insights; and 3) a 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.

References

- 11] 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, Sibo Song, Kai Dang, Peng Wang,
 Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi 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.
- [4] 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, pages 9650–9660, 2021.
- [5] 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 <u>International Conference on</u>
 Machine Learning, pages 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, 2023.
- 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, 2023. URL https://arxiv.org/abs/2307.00300.
- [8] 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, pages 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 <u>ICLR</u>, 2024.
- 338 [11] Alexey Dosovitskiy and Thomas Brox. Generating images with perceptual similarity metrics based on deep networks. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, editors, Advances in Neural Information Processing Systems, volume 29. Curran Associates, Inc., 2016. URL https://proceedings.neurips.cc/paper_files/paper/2016/file/371bce7dc83817b7893bcdeed13799b5-Paper.pdf.
- 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, pages 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, pages 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, 2023. URL https://openreview.net/forum?id=NAQvF08TcyG.
- If all 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, 2023.
- Inhwa Han, Serin Yang, Taesung Kwon, and Jong Chul Ye. Highly personalized text embedding for image manipulation by stable diffusion. arXiv preprint arXiv:2303.08767, 2023.
- Shaozhe Hao, Kai Han, Shihao Zhao, and Kwan-Yee K Wong. Vico: Plug-and-play visual condition for personalized text-to-image generation. arXiv preprint arXiv:2306.00971, 2023.
- Junjie He, Yuxiang Tuo, Binghui Chen, Chongyang Zhong, Yifeng Geng, and Liefeng Bo. Anystory: Towards unified single and multiple subject personalization in text-to-image generation, 2025. URL https://arxiv.org/abs/2501.09503.
- [20] Hexiang Hu, Kelvin C. K. Chan, Yu-Chuan Su, Wenhu Chen, Yandong Li, Kihyuk Sohn, Yang
 Zhao, Xue Ben, Boqing Gong, William Cohen, Ming-Wei Chang, and Xuhui Jia. Instructimagen: Image generation with multi-modal instruction, 2024. URL https://arxiv.org/abs/2401.01952.
- 271 [21] Hexiang Hu, Kelvin CK Chan, Yu-Chuan Su, Wenhu Chen, Yandong Li, Kihyuk Sohn, Yang Zhao, Xue Ben, Boqing Gong, William Cohen, et al. Instruct-imagen: Image generation with multi-modal instruction. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 4754–4763, 2024.
- 375 [22] Miao Hua, Jiawei Liu, Fei Ding, Wei Liu, Jie Wu, and Qian He. Dreamtuner: Single image is enough for subject-driven generation, 2023. URL https://arxiv.org/abs/2312.13691.
- Linyan Huang, Haonan Lin, Yanning Zhou, and Kaiwen Xiao. Flexip: Dynamic control of preservation and personality for customized image generation, 2025. URL https://arxiv.org/abs/2504.07405.
- Zhipeng Huang, Shaobin Zhuang, Canmiao Fu, Binxin Yang, Ying Zhang, Chong Sun, Zhizheng
 Zhang, Yali Wang, Chen Li, and Zheng-Jun Zha. Wegen: A unified model for interactive multimodal generation as we chat, 2025. URL https://arxiv.org/abs/2503.01115.
- Minguk Kang, Jun-Yan Zhu, Richard Zhang, Jaesik Park, Eli Shechtman, Sylvain Paris, and
 Taesung Park. Scaling up gans for text-to-image synthesis. In Proceedings of the IEEE/CVF
 conference on computer vision and pattern recognition, pages 10124–10134, 2023.
- [26] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson,
 Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In
 Proceedings of the IEEE/CVF international conference on computer vision, pages 4015–4026,
 2023.
- Yuval Kirstain, Adam Polyak, Uriel Singer, Shahbuland Matiana, Joe Penna, and Omer Levy.
 Pick-a-pic: An open dataset of user preferences for text-to-image generation. <u>Advances in</u>
 Neural Information Processing Systems, 36:36652–36663, 2023.
- Nupur Kumari, Bingliang Zhang, Richard Zhang, Eli Shechtman, and Jun-Yan Zhu. Multiconcept customization of text-to-image diffusion. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 1931–1941, 2023.
- [29] Duong H. Le, Tuan Pham, Sangho Lee, Christopher Clark, Aniruddha Kembhavi, Stephan
 Mandt, Ranjay Krishna, and Jiasen Lu. One diffusion to generate them all, 2024. URL
 https://arxiv.org/abs/2411.16318.
- Dongxu Li, Junnan Li, and Steven Hoi. Blip-diffusion: Pre-trained subject representation for controllable text-to-image generation and editing. Advances in Neural Information Processing Systems, 36:30146–30166, 2023.

- 402 [31] Zhen Li, Mingdeng Cao, Xintao Wang, Zhongang Qi, Ming-Ming Cheng, and Ying Shan.
 403 Photomaker: Customizing realistic human photos via stacked id embedding, 2023. URL
 404 https://arxiv.org/abs/2312.04461.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr
 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In Computer
 vision–ECCV 2014: 13th European conference, zurich, Switzerland, September 6-12, 2014,
 proceedings, part v 13, pages 740–755. Springer, 2014.
- [33] Zheyuan Liu, Cristian Rodriguez-Opazo, Damien Teney, and Stephen Gould. Image retrieval on
 real-life images with pre-trained vision-and-language models. In Proceedings of the IEEE/CVF
 International Conference on Computer Vision, pages 2125–2134, 2021.
- Zhiheng Liu, Yifei Zhang, Yujun Shen, Kecheng Zheng, Kai Zhu, Ruili Feng, Yu Liu, Deli
 Zhao, Jingren Zhou, and Yang Cao. Cones 2: Customizable image synthesis with multiple
 subjects, 2023. URL https://arxiv.org/abs/2305.19327.
- [35] Jian Ma, Junhao Liang, Chen Chen, and Haonan Lu. Subject-diffusion:open domain personal ized text-to-image generation without test-time fine-tuning, 2024. URL https://arxiv.org/abs/2307.11410.
- Zhendong Mao, Mengqi Huang, Fei Ding, Mingcong Liu, Qian He, and Yongdong Zhang.
 Realcustom++: Representing images as real-word for real-time customization. arXiv preprint arXiv:2408.09744, 2024.
- 421 [37] OpenAI. Introducing gpt-4o and more tools to chatgpt free users, 2024. URL https://openai. 422 com/index/gpt-4o-and-more-tools-to-chatgpt-free/. Accessed: 2024-06-15.
- 423 [38] Or Patashnik, Rinon Gal, Daniil Ostashev, Sergey Tulyakov, Kfir Aberman, and Daniel Cohen-424 Or. Nested attention: Semantic-aware attention values for concept personalization, 2025. URL 425 https://arxiv.org/abs/2501.01407.
- [39] Maitreya Patel, Sangmin Jung, Chitta Baral, and Yezhou Yang. λ-eclipse: Multi-concept
 personalized text-to-image diffusion models by leveraging clip latent space. arXiv preprint
 arXiv:2402.05195, 2024.
- 429 [40] William Peebles and Saining Xie. Scalable diffusion models with transformers, 2023. URL https://arxiv.org/abs/2212.09748.
- [41] Yuang Peng, Yuxin Cui, Haomiao Tang, Zekun Qi, Runpei Dong, Jing Bai, Chunrui Han,
 Zheng Ge, Xiangyu Zhang, and Shu-Tao Xia. Dreambench++: A human-aligned benchmark
 for personalized image generation. In <u>The Thirteenth International Conference on Learning</u>
 Representations, 2025. URL https://openreview.net/forum?id=4GS0ESJrk6.
- 435 [42] pin. https://www.pinterest.com/. https://www.pinterest.com/.
- [43] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe
 Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image
 synthesis, 2023. URL https://arxiv.org/abs/2307.01952.
- Ekta Prashnani, Hong Cai, Yasamin Mostofi, and Pradeep Sen. Pieapp: Perceptual image-error assessment through pairwise preference. In 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 1808–1817, 2018. doi: 10.1109/CVPR.2018.00194.
- [45] Zeju Qiu, Weiyang Liu, Haiwen Feng, Yuxuan Xue, Yao Feng, Zhen Liu, Dan Zhang, Adrian
 Weller, and Bernhard Schölkopf. Controlling text-to-image diffusion by orthogonal finetuning.
 Advances in Neural Information Processing Systems, 36:79320–79362, 2023.
- [46] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
 models from natural language supervision. In <u>International conference on machine learning</u>,
 pages 8748–8763. PmLR, 2021.

- 449 [47] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-450 resolution image synthesis with latent diffusion models. In <u>Proceedings of the IEEE/CVF</u> 451 conference on computer vision and pattern recognition, pages 10684–10695, 2022.
- 452 [48] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-453 resolution image synthesis with latent diffusion models, 2022. URL https://arxiv.org/abs/ 454 2112.10752.
- [49] Ciara Rowles, Shimon Vainer, Dante De Nigris, Slava Elizarov, Konstantin Kutsy, and Simon
 Donné. Ipadapter-instruct: Resolving ambiguity in image-based conditioning using instruct
 prompts, 2024. URL https://arxiv.org/abs/2408.03209.
- In Staniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman.
 Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In
 Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages
 22500–22510, 2023.
- Lolitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al.
 Photorealistic text-to-image diffusion models with deep language understanding. Advances in neural information processing systems, 35:36479–36494, 2022.
- Instantbooth: Personalized text-to-image
 generation without test-time finetuning, 2023. URL https://arxiv.org/abs/2304.03411.
- Quan Sun, Yufeng Cui, Xiaosong Zhang, Fan Zhang, Qiying Yu, Yueze Wang, Yongming Rao, Jingjing Liu, Tiejun Huang, and Xinlong Wang. Generative multimodal models are incontext learners. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 14398–14409, 2024.
- 472 [54] Shangkun Sun, Bowen Qu, Xiaoyu Liang, Songlin Fan, and Wei Gao. Ie-bench: Advancing
 473 the measurement of text-driven image editing for human perception alignment. arXiv preprint
 474 arXiv:2501.09927, 2025.
- Zhenxiong Tan, Songhua Liu, Xingyi Yang, Qiaochu Xue, and Xinchao Wang. Ominicontrol:
 Minimal and universal control for diffusion transformer. arXiv preprint arXiv:2411.15098,
 2024.
- 478 [56] uns. https://unsplash.com/. https://unsplash.com/.
- 479 [57] Andrey Voynov, Qinghao Chu, Daniel Cohen-Or, and Kfir Aberman. P+: Extended textual conditioning in text-to-image generation, 2023. URL https://arxiv.org/abs/2303.09522.
- [58] Haofan Wang, Matteo Spinelli, Qixun Wang, Xu Bai, Zekui Qin, and Anthony Chen. Instantstyle: Free lunch towards style-preserving in text-to-image generation. <u>arXiv:2404.02733</u>, 2024.
- 484 [59] Xierui Wang, Siming Fu, Qihan Huang, Wanggui He, and Hao Jiang. Ms-diffusion: Multi-485 subject zero-shot image personalization with layout guidance. <u>arXiv preprint arXiv:2406.07209</u>, 486 2024.
- 487 [60] Yibin Wang, Yuhang Zang, Hao Li, Cheng Jin, and Jiaqi Wang. Unified reward model for multimodal understanding and generation. arXiv preprint arXiv:2503.05236, 2025.
- [61] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi,
 Quoc V. Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language
 models. In Advances in Neural Information Processing Systems (NeurIPS), 2022.
- Yuxiang Wei, Yabo Zhang, Zhilong Ji, Jinfeng Bai, Lei Zhang, and Wangmeng Zuo. Elite:
 Encoding visual concepts into textual embeddings for customized text-to-image generation.
 In Proceedings of the IEEE/CVF International Conference on Computer Vision, pages 15943–15953, 2023.

- 496 [63] Wikipedia. https://en.wikipedia.org/. https://en.wikipedia.org/. [Online; accessed 11-497 May-2025].
- Shaojin Wu, Mengqi Huang, Wenxu Wu, Yufeng Cheng, Fei Ding, and Qian He. Less-to-more generalization: Unlocking more controllability by in-context generation. <u>arXiv preprint</u> arXiv:2504.02160, 2025.
- 501 [65] Xiaoshi Wu, Keqiang Sun, Feng Zhu, Rui Zhao, and Hongsheng Li. Human preference score:
 502 Better aligning text-to-image models with human preference. In <u>Proceedings of the IEEE/CVF</u>
 503 International Conference on Computer Vision, pages 2096–2105, 2023.
- [66] Shitao Xiao, Yueze Wang, Junjie Zhou, Huaying Yuan, Xingrun Xing, Ruiran Yan, Chaofan
 Li, Shuting Wang, Tiejun Huang, and Zheng Liu. Omnigen: Unified image generation. arXiv
 preprint arXiv:2409.11340, 2024.
- 507 [67] Zhexiao Xiong, Wei Xiong, Jing Shi, He Zhang, Yizhi Song, and Nathan Jacobs. Grounding-508 booth: Grounding text-to-image customization, 2025. URL https://arxiv.org/abs/2409. 509 08520.
- [68] Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao
 Dong. Imagereward: Learning and evaluating human preferences for text-to-image generation.
 Advances in Neural Information Processing Systems, 36:15903–15935, 2023.
- [69] Jiazheng Xu, Yu Huang, Jiale Cheng, Yuanming Yang, Jiajun Xu, Yuan Wang, Wenbo Duan,
 Shen Yang, Qunlin Jin, Shurun Li, Jiayan Teng, Zhuoyi Yang, Wendi Zheng, Xiao Liu, Ming
 Ding, Xiaohan Zhang, Xiaotao Gu, Shiyu Huang, Minlie Huang, Jie Tang, and Yuxiao Dong.
 Visionreward: Fine-grained multi-dimensional human preference learning for image and video
 generation, 2024. URL https://arxiv.org/abs/2412.21059.
- 518 [70] Hu Ye, Jun Zhang, Sibo Liu, Xiao Han, and Wei Yang. Ip-adapter: Text compatible image prompt adapter for text-to-image diffusion models. arXiv preprint arXiv:2308.06721, 2023.
- 520 [71] Yu Zeng, Vishal M. Patel, Haochen Wang, Xun Huang, Ting-Chun Wang, Ming-Yu Liu, and 521 Yogesh Balaji. Jedi: Joint-image diffusion models for finetuning-free personalized text-to-image 522 generation, 2024. URL https://arxiv.org/abs/2407.06187.
- Fig. 172 Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 586–595, 2018. doi: 10.1109/CVPR.2018.00068.
- Yuxuan Zhang, Yiren Song, Jiaming Liu, Rui Wang, Jinpeng Yu, Hao Tang, Huaxia Li, Xu Tang, Yao Hu, Han Pan, et al. Ssr-encoder: Encoding selective subject representation for subject-driven generation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 8069–8078, 2024.

A Details of Hierarchical Category Establishing

531

532

533

534

535

536

537

538

539

543

544

545

546

547

549

550

551

552

553

554

557

558

559

560

561

562 563

564

565

566

567

568

The First-level Category We observed the composition of existing benchmark data. From a more abstract and higher-level perspective 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-40 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 7.

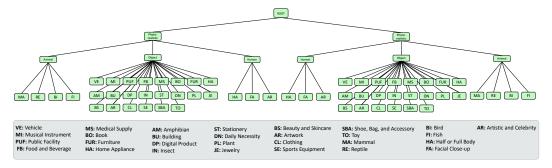


Figure 7: **The hierarchical category system.** We developed a three-level category hierarchy by integrating data from existing large-scale datasets and open-source encyclopedic resources.

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

pen pencil fax machine stapler vehicle car van airplane motorbike saibboat unicycle motorcycle motorbike saibboat unicycle motorbike saibboard picture frame movie (disc) playing cards table cloth telephone laptop computer remote mouse keyboard printer desktop copier remote with toner blush serum emulsion sunscreen bottle spoon clock toothbrush basket pillow power outlet light switch basket player of train boat	fity firefly ant butterfly ladybug locust dragonfly Insect amphibian frog bullfrog toad salamander			Candida	te Category Labels	i			The Third-level Categor
mamphibian frog bullfrog toad salamander Amphibian fish goldfish seahorse shark tilapia Fish Fish bird chicken turkey swallow crow pigeon Soose rooster Bird mammal cat squirrel bear squirrel gard early shale shark Soos and pigeon Soose rooster mammal cat squirrel dog giraffe lion monkey tiger bunny playpus whale aardvark rabbis Sofa sandwark Fish basees street field goal post offee table window door libaseball court baselenble confee table window Sofa table furniture dining table sofa chair couch bed desk table window Sofa table furniture potent plant tree Sunflower Cacuts lavender Plant cookie milk pineapple carrot broccoli banana pineapple carrot broccoli banana cheese cupeake doont dryer firidge refrigerator sheed wheelchair pants Saket shorts Sos shorts Sofa mammal cheekace bracelet shakethall shower some shorts	mamphibian frog bullfrog toad salamander Amphibian fish goldfish seahorse shark tilapia Fish bird chicken turkey swallow crow pigeon Fish bird chicken turkey swallow crow pigeon Fish mammal cat sopired giraffe lion monkey tiger bunny platypus whale sardwark rabbit Fish street fountain fire hydrant stable lamp Fish furniture dining table sofa chair cache window door coffee table window door furniture plotted plant tree sunflower cactus lavender Furniture cookie pineapple carrot broccoli cake pizza soup meat pumpkin cheese cupeake donnt hold ong pinter fan (ceil/floor) printer fax machine copier meklace bracelet ring pendant shorts safe lagset lang sleeve shirt shorts safe lagset lang sleeve shirt shorts shorts safe lang sleeve shirt shorts shorts lagset lang sleeve shirt shorts lagset l	reptile	lizard	dinosaur	turtle	crocodile	chameleon	gecko	Reptile
fish goldrish seaborse shark tilapia seaborse Fish bird chicken duck owd swan goose rooster Bird mammal cat duck swallow crow pigeon croster Bird mammal cat squirrel godg horse sheep cow clephant goat squirrel gigariffe kangaroo rabin swall horse sheep cow clephant street fountain fire hydrant traffic light bus stop sign sign parking meter goal net Public Facility furniture dining table coffee table side table chench cabinet mirror carpet Furniture flower potted plant tree sundiche bacch cabit parker baca parter carpet Furniture cookie milk pancake pasca grape cereal bean <td< td=""><td>fish goldfish seahorse shark tilapia good Fish bird chicken duck owl cow psyan goose rooster Bird mammal cat dog horse sheep cow elephant goat pig squirrel graffe lion monkey tiger bunny street fountain fire hydrant traffic light bus stop sign parking meter goal net Public Facility furniture dining table table coffee table window door chair couch bed desk parking meter goal net Public Facility flower potted chandelier table bench cabinet mirr cape Furniture down milk paacake pasta grape cereal pate pate pate parker pate pate pate pate pumpkin Food and Beverag debender fridge <</td><td>fly</td><td>firefly</td><td>ant</td><td>butterfly</td><td>ladybug</td><td>locust</td><td>dragonfly</td><td>Insect</td></td<>	fish goldfish seahorse shark tilapia good Fish bird chicken duck owl cow psyan goose rooster Bird mammal cat dog horse sheep cow elephant goat pig squirrel graffe lion monkey tiger bunny street fountain fire hydrant traffic light bus stop sign parking meter goal net Public Facility furniture dining table table coffee table window door chair couch bed desk parking meter goal net Public Facility flower potted chandelier table bench cabinet mirr cape Furniture down milk paacake pasta grape cereal pate pate pate parker pate pate pate pate pumpkin Food and Beverag debender fridge <	fly	firefly	ant	butterfly	ladybug	locust	dragonfly	Insect
bird hen turkey swallow crow pigeon pigeon washer squired goat whale art washer goat winder and washer shorts cover in the stethoscope book magazine book magazine book magazine shorts while book magazine fijn plop book magazine textbook motorbycle motorcycle motorbike skateball motorcycle motorcycle motorcycle motorbike skateball motorcycle moto	bird hen turkey swallow crow pigeon pigeon rooster birder goat pig kangaroo whale aratheria baschall court baskethall court turkey window door chandelier table window door chandelier table window door chandelier table blender birdered sandwich cake donut both dog bannana orange strawberry apple pumpkin cheese cupcake donut both dog bannana orange strawberry apple pumpkin cheese cupcake donut both dog bannana orange strawberry apple pumpkin cheese cupcake donut both dog bannana orange strawberry apple pumpkin cheese cupcake donut both dog bannana orange strawberry apple pumpkin cheese cupcake donut both dog bannana orange strawberry apple pumpkin shorts swallow blender blender banker title lamp gate soup meat pumpkin shorts swallow blender bench cheese cupcake donut both dog bannana orange strawberry apple pumpkin shorts swallow to the dog bannana orange strawberry apple pumpkin shorts swallow to the dog bannana orange strawberry apple pumpkin shorts swallow to the dog bannana orange strawberry apple pumpkin shorts swallow to the dog bannana orange strawberry apple pumpkin shorts swallow to the dog bannana orange strawberry apple pumpkin shorts swallow to the dog bannana orange strawberry apple pumpkin short seed to the dog bannana orange strawberry apple pumpkin short seed to the dog bannana orange strawberry apple pumpkin short seed bracelet ring pendant brooch anklet washer (blook short seed short seed to the super hero costume sock swall short seed short see	amphibian	frog	bullfrog	toad	salamander			Amphibian
mammal bear squirrel giraffe giraffe ribinon monkey tiger blumny platypus solution whale sacret squirrel goal metabare squirrel goal sardvark metabaseball court busseps since to baseball court busseps since the table and court beat table and court table and court beat table and court beat table and court table and court beat table and court tabl	mammal bear squirrel giraffe shangaroo marda sada saa saa saa saa saa saa saa saa	fish	goldfish	seahorse	shark	tilapia			Fish
mammal cat dog many and cat gord per	mammal cat dog graffe goat whale arriver goal reference bear squared goat whale arriver field goal post soccer net baskethall court baskethal basketha						goose	rooster	Bird
bear goat whale squirrel goat whale climation or substitution of the point of	bear goat whale squirrel pig aardvark giraffe kangaroo arabbit lion core zebra monkey deer mouse tigger platypus bunny platypus Mammal mammal platypus field goal post street field goal post street window fountain soccernet basketball cout but stop sign chair couch cabinet cabinet beach bed desk cabinet beach cabinet cabinet beach cabinet beach door goal net platypus Public Facility flower potted plant beach door tree sunflower cactus lavender		turkey	swallow	crow	pigeon) Jiiu
goat whale pig aardwark kangaroo rabbit rhinoceros mouse deer mouse hippo platyfus platyfus Mammal plate field goal post soccer net basketball court bus stop sign sign parking meter goal net goal	goat whale 'pig ardvark kangaroo rabbit rininoceros zebra deer bit pour la bit pour la bit pour la basketball court kangaroo basketball court rininoceros zebra deer bit pour la								
whale aardvark rabbit zebra mouse " " Public Facility field goal post soccer net fold goal post fountain fable soccer net basketball court basketball c	whale aardvark rabbit zebra mouse " " " " " " " " " " " " " " " " " " "								Mammal
furniture dining table sofa table don't don't claimly shown to be table window door door door door claimly shown to door door door door door claimly shown door door door door claimly shown door door door claimly shown door door door claimly shown door door door door door door door doo	furniture dining table table coffee table co							I J I	
table window coffee table door side table chandelier bench table lamp gate cabinet gate mirror gate carpet Furniture flower potted plant tree sunflower cactus lavender Plant cookie pineapple bread codes milk pancake carrot sandwich cake piread cheese pasta pancake pirza soup bacon cereal strawberry pumpkin carbot tomato bean apple pumpkin pipiza soup pumpkin carbot tomato Food and Beverage dryer blender friidge particular from (ceil/floor) microwave pirza soup bacon oven coster copier washer Home Appliance necklace bracelet ring pendant brooch anklet Jewelry wheelchair gauze crutch stethoscope syringe washer Home Appliance book magazine textbook dictionary biography birt Clothing bat shorts skis snowboard tennis racket baketball hoop baseball glove soccer ball Sports Equipment flip flop handbag glove shoe backpack washer S	table window of door chandelier table lamp gate mirror carpet Furniture window of door chandelier table lamp gate gate mirror carpet furniture gate potted plant tree sunflower cactus lavender plant tree sunflower cactus lavender plant tree carot broccoli banana orange strawberry pumpkin carbo bread sandwich cake care broccoli banana orange strawberry pumpkin carbo broccoli bender carot donut hot dog bacon egg trawberry pumpkin tomato pumpkin tomato donut hot dog bacon egg trawberry pumpkin tomato pumpkin tomato donut hot dog bacon egg trawberry pumpkin tomato pumpkin tomato donut hot dog bacon egg trawberry pumpkin tomato pumpkin tomato donut pumpkin tomato pumpkin tomato donut donut hot dog bacon egg trawberry pumpkin tomato papic pumpkin tomato donut donut					sign	parking meter	goal net	Public Facility
window door chandelier table lamp gate " Plant flower potted plant tree sunflower cactus lavender Plant cookie pineapple bried plant pineapple bread cheese milk pineapple carrot broccoli sandwich cake pasta pasta pizza banana papple pizza soup care pumpkin apple apple pumpkin apple pumpkin apple pumpkin apple pumpkin apple apple apple apple apple pumpkin apple pumpkin apple	Bower	furniture	dining table	sofa	chair	couch	bed	desk	<u> </u>
flower potted plant tree sunflower cactus lavender footbell carrot carbo standwich carbo sandwich carbo sandwich cake broccoli banana orange strawberry strawberry sandwich cake broccoli banana orange strawberry strawberry strawberry strawberry strawberry strawberry strawberry should be carrot cake broccoli banana orange strawberry strawberry pumpkin bean orange grape orange strawberry strawberry pumpkin bean orange strawberry strawberry pumpkin bean papte pumpkin tomato strawberrage football banana orange strawberry pumpkin tomato strawberrage football should be strawberrage football should be strawberrage for the steep soup meat to toaster fax machine copier fax machine copier fax machine copier sock in the steep soup strawberrage fax machine steep shirt short sleeve shirt short short short shall basketball football tennis racket tennis net hoop baseball glove soccer ball sports Equipment short shall basketball shop basket	flower potted plant tree sunflower cactus lavender Plant cookie pineapple carrot broccoli banana orange strawberry apple plant sandwich cake sandwich cheese cupcake donut hot dog bacon egg trawberry apple pumpkin tomato dryer fridge hair drier fan (ceil/floor) printer fax machine copier necklace bracelet ring pendant broccol anklet Jewelry wheelchair gauze crutch stethoscope syringe wheelchair scarf lie super hero costume sock shorts scarf some book magazine textbook dictionary biography book magazine textbook dictionary biography flip flop handbag glove shoe backpack Shoe, Bag, and Access pen pencil fax machine stapler vehicle car airplane hotorcycle motorcycle motorcycle motorcycle motorcycle motorcycle motorcycle motorcycle from flute violin for skiete laptop from the flute violin for sunscreen from the flute violin for sunscreen from the flute violin for sunscreen facel to show a sunscreen facel to plate from the flute cup bowl teapot chopping board ladder basket basket ball book for king for k						mirror	carpet	Furniture
cookie pineapple carrot broccoli brand carbot carbot broccoli brand cheese cupcake donut hot dog brand cheese cupcake donut hot dog bacon egg strawberry pumpkin cheese cupcake donut hot dog bacon egg strawberry pumpkin cheese cupcake donut hot dog bacon egg broth tomato cupcake donut hot dog bacon egg brand tomato cupcake bat skis shorts ega brand tie super hero costume sock bat skis snowboard tennis racket tennis net hoop baseball glove soccer ball basketball basketball basketball football tennis racket tennis net hoop baseball glove soccer ball sports Equipment egg brand tennis racket tennis net hoop backpack Shoe, Bag, and Accesso pen pencil fax machine stapler Stationery ehicle sailboat raft airplane helicopter hot air balloon rocket bicycle element with the element process table cloth Shoe, Bag, and Accesso pen brand truck backpack Shoe, Bag, and Accesso pen brand truck backpack Shoe, Bag, and Accesso pen element emotorbyle skateboard hot air balloon rocket bicycle element elephone motorcycle motorcycle motorcycle motorcycle fat table cloth Shoe, Bag, and Accesso pen element elephone element elephone elapto computer tablet ipad iphone cell phone radio Digital Product kite toy cars toy legos robot doll Musical Instrument telephone clock toothbrush basket board towal candle balloon power outlet balloon power outlet blight switch Daily Necessity	cookie pineapple bread carrot bread cheese milk carrot sandwich cheese pancake broccoli cake donut pasta banana pizza grape orange soup bacon cereal egg bean apple pumpkin tomato Food and Beverage pumpkin tomato dryer blender fridge hair drier refrigerator fan (ceil/floor) microwave printer oven fax machine toaster copier washer Home Appliance necklace bracelet ring pendant brooch anklet Jewelry wheelchair gauze crutch stethoscope syringe washer Home Appliance pants jacket long sleeve shirt short sleeve shirt pajamas underpants shirt Clothing book magazine textbook dictionary biography soccer ball Sports Equipmen flip flop handbag glove shoe backatball hoop baseball glove soccer ball Sports Equipmen flip flop handbag glove shoe backatball hoop baseball glove soccer ball Sports Equipmen vehicle car van truc	window	door	chandelier	table lamp	gate			
pineapple bread sandwich cake pizza borocoli cheese cupcake donut hot dog bacon egg strawberry apple pumpkin tomato cupcake donut hot dog bacon egg washer tomato food and Beverage blender blender hair drier fan (ceil/floor) microwave printer fax machine copier sowe printer fax machine copier sowe printer fax machine copier washer looked bracelet ring pendant brooch anklet Jewelry wheelchair gauze crutch stethoscope syringe Medical Supply pants jacket long sleeve shirt short sleeve shirt short sleeve shirt short sleeve shirt sock scarf tite super hero costume sock underpants shorts book magazine textbook dictionary biography Book bat skis snowboard football tennis racket tennis racket tennis racket tennis reacket tennis printer shorts ball basketball football football football raft airplane airplane helicopter skateboard motorcycle motorcycle motorcycle motorbike skateboard skateboard bouse building roof bridge church Building picture frame movie (disc) playing cards table cloth Fight tablet to y cars toy legos robot doll hair drum flute violin toward ladder bowl teapot copier add ladder bowl bowl teapot copier and in phone cell phone radio Digital Product bowl bowl teapot copier add loll bowl bowl toward cook chopping board ladder bowl vase towel candle basket balloon power outlet light switch balloon chopping board ladder basket pillow power outlet light switch balloon box chopping board ladder basket pillow power outlet light switch balloon box chopping board ladder basket pillow power outlet light switch balloon box ochoping board ladder basket pillow power outlet light switch balloon box ochoping board ladder basket pillow power outlet light switch balloon box ochoping board ladder basket pillow power outlet light switch balloon box ochoping board ladder basket pillow power outlet light switch balloon basket balloon power outlet light switch balloon basket balloon power outlet light switch balloon basket balloon power outlet light switch balloon by the power outlet light switch balloon by the power ou	pineapple bread carrot sandwich cake pizza broad cake pizza soup soup meat pumpkin cheese cupcake donut hot dog bacon egg tomato food and Beverage dryer blender fair drier fan (ceil/floor) microwave printer fax machine copier washer Home Appliance fax machine copier fax (ceil/floor) microwave printer fax machine copier fax washer fax machine copier fax washer Jewelchair gauze crutch stethoscope syringe Medical Supply pants jacket shorts searf tie super hero costume sock football shaketball book magazine textbook dictionary biography Sports ball basketball football tennis racket sports ball basketball football tennis net backpack Shoe, Bag, and Access pen pencil fax machine sailboat unicycle motorbike motorbike motorbike motorbike motorbike motorbike skateboard motoreycle motorbike skateboard shaketball pipting cards table cloth football lastrument guitar drum flute violin Musical Instrument telephone mouse keyboard printer desktop copier radio Digital Product kite toy cars toy legos robot doll bould bat back table coup bowl teapot fork kinife spoon clock toothbrush vase towel light switch light switch bails on power outlet light switch baily now power outlet light switch light switch printer appeal	flower	potted plant	tree	sunflower	cactus	lavender		Plant
bried cheese cupcake donut hot dog bacon egg pumpkin cupcake donut hot dog bacon egg dryer fridge cupcake donut hot dog bacon egg pumpkin formato fram cipcial fram drier fax machine copier printer fax machine copier pen pencil fax machine salebat move for gridge and per pencil fax machine salebat move from the salebat move from	cheese cupcake donut hot dog bacon egg pumpkin tomato dout tomato hot dog bacon egg pumpkin tomato dout tomato dout tomato hot dog bacon egg pumpkin tomato dout tomato dout tomato hot dog bacon egg pumpkin tomato dout tomato dout tomato printer fan (ceil/floor) printer fan copier printer fan copier printer fan copier printer gauze crutch stethoscope syringe Medical Supply pants jacket scarf tie super hero costume sock will be skis scarf tie super hero costume sock will be skis sonwboard tennis net hoop baseball glove soccer ball basketball football tennis net hoop baseball glove soccer ball hoop baseball glove soccer ball basketball football tennis net hoop will be sketball football football tennis net hoop will be sketball football								
cheese cupcake donut hot dog bacon egg tomato Pood and Beverage dyer blender fridge hair drier fan (ceil/floor) printer fan washer copier washer Home Appliance fan (ceil/floor) printer fan washer copier washer Home Appliance over the fan washer copier washer washer copier washer Home Appliance fan washer copier washer washer copier washer Home Appliance fan washer copier washer Home Appliance washer copier washer washer copier washer washer copier washer w	cheese cupcake donut hot dog bacon egg tomato Pood and Beverage dryer blender firidge refrigerator fan (ceil/floor) microwave printer fax machine copier washer hair drier fax (ceil/floor) microwave printer fax machine copier washer lowed fax machine shorts a scarf lie super hero costume sook underpants shirt shorts scarf lie super hero costume sook underpants shirt lie super hero costume sook underpants shirt lowed fax machine lowed for tennis racket tennis racket tennis racket tennis net loop handbag glove shoe backpack Shoe, Bag, and Access pen pencil fax machine stapler Stationery vehicle car avan anipplane unicycle motorcycle motorbike skateboard lowed for bridge church louse building roof bridge church louse building roof bridge church louse laptop computer frame movie (disc) playing cards table cloth louse laptop computer mouse keyboard printer desktop copier radio bigital Product kite toy cars toy legos robot doll hair brush toner blush serum emulsion sunscreen Beauty and Skinca chopping board ladder basket pillow power outlet light switch light switch lought Necessity								
dryer blender fairder fan (ceil/floor) microwave printer fax machine copier washer lecklace bracelet ring pendant brooch anklet Jewelry wheelchair gauze crutch stethoscope syringe Medical Supply pants jacket scarf tie super hero costume sook magazine textbook dictionary biography book magazine textbook dictionary biography basketball basketball football tennis racket tennis net thoop handbag glove shoe backpack Shoe, Bag, and Accesso pen pencil fax machine stapler Stationery which sailboat unicycle motorcycle motorbike skateboard suited instrument guitar drum flute violin Musical Instrument telephone remote mouse keyboard printer desktop copier radio Digital Product kite toy cars toy legos robot doll bowl chopping board ladder basket basket pillow power outlet light switch baily Necessity towel candle basket balloon power outlet light switch baily Necessity	dryer blender fridge hair drier fan (ceil/floor) microwave printer fax machine copier washer later fax (ceil/floor) printer fax machine copier washer later fax machine copier mecklace bracelet ring pendant brooch anklet Jewelry meelchair gauze crutch stethoscope syringe Medical Supply pants shorts scarf lie super hero costume sock shorts scarf lie super hero costume sock shorts skis shorts skis snowboard sports ball basketball football tennis racket tennis racket tennis racket tennis racket loop handbag glove shoe backpack Shoe, Bag, and Access pen pencil fax machine stapler skateboard motorcycle motorcycle motorbike skateboard loop building picture frame movie (disc) playing cards table cloth mouse keyboard printer desktop copier remote mouse keyboard pour letges round for the spons box chopping board ladder basket pillow power outlet light switch light switch loop power outlet light switch loader ladder basket pillow power outlet light switch loader light switch loader ladder basket pillow power outlet light switch loader light switch loader light switch loader ladder basket pillow power outlet light switch loader loader loader loader loader loader loader ladder basket pillow power outlet light switch loader light switch loader ladder basket pillow power outlet light switch loader ladder ladder ladder loader printer loader ladder ladder ladder loader ladder loader ladder loader ladder ladder ladder loader ladder loader ladder loader ladder ladder ladder ladder ladder ladder ladder ladder ladder la								Food and Beverage
Definder hair drier fan (ceil/floor) printer fax machine copier Frome Appliance	Definder hair drier fan (ceil/floor) printer fax machine copier Profile Apphainter		*						
mecklace bracelet ring pendant brooch anklet Jewelry wheelchair gauze crutch stethoscope syringe Medical Supply pants jacket scarf lite super hero costume sock underpants shirt Short sleeve shirt sock which sock which short shor	mecklace bracelet ring pendant brooch anklet Jewelry wheelchair gauze crutch stethoscope syringe Medical Supply pants jacket scarf lie super hero costume shorts searf scarf tie super hero costume sock underpants shirt shorts searf shorts searf tie super hero costume sock underpants shirt Shorts searf scarf lie super hero costume sock underpants shirt Shorts searf scarf super hero costume sock underpants shirt Shorts searf shorts searf scarf lie super hero costume sock underpants shirt Shorts searf scarf super hero costume sock underpants shirt Shorts searf short sleeve shirt shorts searf scarf super hero costume sock underpants shirt Shorts searf scarf short sleeve shirt shorts searf scarf shorts searf sock lie spatial hoop baseball glove soccer ball sports Equipmen shorts hoop hoop saseball glove soccer ball sports Equipmen shorts hoop hoop saseball glove soccer ball sports Equipmen shorts hoop hoop saseball glove soccer ball sports Equipmen shorts hoop hoop saseball glove soccer ball sports Equipmen shorts hoop hoop saseball glove soccer ball sports Equipmen shorts hoop hoop saseball glove soccer ball sports Equipmen shorts hoop hoop haseball glove soccer ball sports Equipmen shorts hoop hoop saseball glove soccer ball sports Equipmen shorts hoop hoop haseball glove soccer ball sports Equipmen shorts hoop hoop haseball glove soccer ball sports Equipmen shorts hoop hoop haseball glove soccer ball shorts hoop hoop haseball glove soccer ball sports Equipmen shorts hoop hoop haseball glove soccer ball shoop hoop haseball glove soccer ball shorts hoop hoop haseball glove soccer bal							washer	Home Appliance
wheelchair gauze crutch stethoscope syringe Medical Supply pants shorts jacket shorts long sleeve shirt tie shorts sock underpants shirt Clothing book magazine textbook dictionary biography Book bat sports ball skis showboard basketball snowboard football tennis racket tennis net tennis racket tennis net hoop basketball hoop baseball glove soccer ball hoop Sports Equipment flip flop handbag glove shoe backpack Shoe, Bag, and Accesso pen pencil fax machine stapler Stationery vehicle sailboat unicycle car raft airplane motorbike skateboard bus train hot air balloon rocket biocycle house building roof bridge church Building picture frame movie (disc) playing cards table cloth was table cloth Musical Instrument telephone remote laptop computer sevente tablet ipad desktop iphone coll phone radio cell phone radio Digital Product	wheelchair gauze crutch stethoscope syringe Medical Supply pants shorts jacket scarf long sleeve shirt tie short super hero costume pajamas sock underpants shirt Clothing book magazine textbook dictionary biography Book bat sports ball skis shetball snowboard football tennis racket tennis net basketball hoop baseball glove soccer ball hoop Sports Equipmen flip flop handbag glove shoe backpack Shoe, Bag, and Acces pen pencil fax machine stapler Stationery vehicle sailboat raft unicycle raft motorbike skateboard bus train hot air balloon rocket boat bicycle Vehicle house building roof bridge church Building Artwork picture frame movie (disc) playing cards table cloth with table cloth Musical Instrument telephone remote laptop computer keyboard printer desktop copier call phone coll phone radio Digital Prod	necklace	bracelet		-	brooch	•		Jewelry
pants shorts scarf long sleeve shirt super hero costume sock underpants shirt Short scarf tite super hero costume sock underpants shirt Short scarf shorts scarf tite super hero costume sock underpants shirt Clothing book magazine textbook dictionary biography Book bat skis snowboard basketball football tennis racket tennis net basketball hoop baseball glove soccer ball sports Equipment flip flop handbag glove shoe backpack Shoe, Bag, and Accessor pen pencil fax machine stapler Stationery vehicle car van airplane helicopter skateboard hot air balloon rocket bicycle unicycle motorcycle motorbike skateboard house building roof bridge church Building picture frame movie (disc) playing cards table cloth Suscial instrument guitar drum flute violin Musical Instrument telephone laptop computer tablet ipad desktop copier radio Digital Product kite toy cars toy legos robot doll hair brush toner blush serum emulsion sunscreen Beauty and Skincare bottle spoon chopping board ladder basket pillow power outlet light switch Daily Necessity	pants shorts								
shorts scarf tie super hero costume sock Clothing book magazine textbook dictionary biography Book bat sports ball skis snowboard football tennis racket tennis net tennis net hoop basketball hoop baseball glove soccer ball hoop Sports Equipment flip flop handbag glove shoe backpack Shoe, Bag, and Accesso pen pencil fax machine stapler Stationery vehicle car van truck bus bus boat sailboat raft airplane helicopter hot air balloon rocket bicycle Vehicle sailboat raft motorcycle motorbike skateboard church Building bouse building roof bridge church Artwork musical instrument guitar drum flute violin Musical Instrument telephone remote laptop mouse keyboard printer desktop <td< td=""><td>shorts scarf tie super hero costume sock Clothing book magazine textbook dictionary biography Book bat skis snowboard football tennis racket sports ball basketball basketball football tennis net hoop baseball glove soccer ball sports Equipmen flip flop handbag glove shoe backpack Shoe, Bag, and Acces pen pencil fax machine stapler vehicle car van airplane motorcycle motorbike skateboard house building roof bridge church bouse building roof bridge church Building picture frame movie (disc) playing cards table cloth rusical instrument guitar drum flute violin Musical Instrument telephone mouse keyboard printer desktop copier radio Digital Product kite toy cars toy legos robot doll hair brush toner blush serum emulsion sunscreen Beauty and Skinca bottle spoon chopping board ladder basket pillow power outlet light switch Daily Necessity</td><td></td><td></td><td></td><td></td><td></td><td>undarmente</td><td>abiet</td><td>1</td></td<>	shorts scarf tie super hero costume sock Clothing book magazine textbook dictionary biography Book bat skis snowboard football tennis racket sports ball basketball basketball football tennis net hoop baseball glove soccer ball sports Equipmen flip flop handbag glove shoe backpack Shoe, Bag, and Acces pen pencil fax machine stapler vehicle car van airplane motorcycle motorbike skateboard house building roof bridge church bouse building roof bridge church Building picture frame movie (disc) playing cards table cloth rusical instrument guitar drum flute violin Musical Instrument telephone mouse keyboard printer desktop copier radio Digital Product kite toy cars toy legos robot doll hair brush toner blush serum emulsion sunscreen Beauty and Skinca bottle spoon chopping board ladder basket pillow power outlet light switch Daily Necessity						undarmente	abiet	1
bat skis sports ball basketball football tennis racket sports ball basketball basketball football tennis net hoop baseball glove soccer ball sports Equipment flip flop handbag glove shoe backpack Shoe, Bag, and Accesso pen pencil fax machine stapler Stationery vehicle car airplane airplane motorcycle motorbike skateboard hot air balloon rocket bicycle who air balloon picture frame movie (disc) playing cards table cloth Artwork who are alaptop computer tablet balloon rocket bicycle who air balloon pinter who are alaptop computer tablet balloon rocket bicycle who are alaptop computer tablet balloon pinter desktop copier radio big plate cup bowl teapot fork knife balloon box chopping board ladder basket pillow power outlet light switch balloon power ou	bat skis snowboard football tennis racket sports ball basketball boop baseball glove soccer ball soccer ball basketball football tennis net basketball hoop baseball glove soccer ball sports Equipmen baseball glove soccer ball sports Equipmen baseball glove soccer ball basketball football tennis net baseball glove soccer ball sports Equipmen backgrows and soccer ball sports Equipmen backgrows and soccer ball sports Equipmen backgrows and scale glove soccer ball sports equipmen to scale glove soccer ball sports and scale glove soccer ball sports glove soccer backgrows glove soccer ball						underpants	SIIII	Clothing
sports ball basketball football tennis net hoop Sports Equipment flip flop handbag glove shoe backpack Shoe, Bag, and Accesso pen pencil fax machine stapler Stationery vehicle car van truck helicopter sailboat unicycle motorbike sailboar unicycle motorbike sailboar unicycle motorbike sateboard house building roof bridge church Sulliang picture frame movie (disc) playing cards table cloth Sulliang sulliar drum flute violin Musical Instrument telephone laptop computer tablet ipad iphone cell phone remote mouse keyboard printer desktop copier radio Digital Product kite toy cars toy legos robot doll hair brush toner blush serum emulsion sunscreen Beauty and Skincare bottle spoon clock toothbrush basket pillow power outlet light switch Daily Necessity	sports ball basketball football tennis net hoop Sports Equipmen flip flop handbag glove shoe backpack Shoe, Bag, and Acces pen pencil fax machine stapler vehicle car van airplane motorcycle motorbike helicopter skateboard house building roof bridge church Suitable cloth house building roof bridge church Suitable cloth fusical instrument guitar drum flute violin Musical Instrumen telephone laptop computer mouse keyboard printer desktop copier radio Digital Product kite toy cars toy legos robot doll hair brush toner blush serum emulsion sunscreen Beauty and Skinca bottle spoon clock chopping board ladder basket pillow power outlet light switch Daily Necessity	book	magazine	textbook	dictionary	biography			Book
pen pencil fax machine stapler vehicle car van airplane motorbike skateboard house building roof bridge church picture frame movie (disc) playing cards table cloth telephone laptop computer remote mouse keyboard telephone laptop computer tablet printer telephone laptop computer tablet biekstop copier radio kite toy cars toy legos robot doll hair brush toner blush serum emulsion sunscreen bottle spoon chopping board ladder basket pillow power outlet light switch Stationery Stationer	pen pencil fax machine stapler vehicle car van airplane helicopter skateboard house building roof bridge church picture frame movie (disc) playing cards table cloth remote laptop computer remote mouse keyboard printer telephone laptop computer remote mouse keyboard bair bair bullon hair brush toner blush serum emulsion sunscreen by stationery Desire frain boat bicycle Stationery Stationery Stationery Stationery Stationery Stationery Stationery Stationery Stationery Desire frain boat bicycle Stationery Stationery Stationery Stationery Desire frain boat bicycle Stationery Stationery Stationery Stationery Desire frain boat bicycle Stationery Toket bicycle Vehicle Stationery Stationery Stationery Desire frain boat bicycle Stationery Toket bicycle Vehicle Stationery Toket bicycle Stationery Toket bicycle Vehicle Stationery Toket bicycle Vehicle Stationery Toket bicycle Stationery T						baseball glove	soccer ball	Sports Equipment
vehicle sailboat unicycle car raft motorcycle van airplane motorbike truck helicopter skateboard hot air balloon train rocket boat bicycle Vehicle house building roof bridge church Image: church with air balloon Building picture frame movie (disc) playing cards table cloth Image: church with air balloon Image: church with air balloon with air balloon Musical Instrument usical instrument guitar drum flute violin Musical Instrument telephone remote laptop computer keyboard tablet printer desktop copier copier radio Digital Product kite toy cars toy legos robot doll hair brush toner blush serum emulsion sunscreen Beauty and Skincare bottle spoon clock chopping board toothbrush ladder vase towel candle balloon power outlet light switch Daily Necessity	vehicle sailboat raft airplane motorbike skateboard skateboard sailboat raft motorbike skateboard s	flip flop	handbag	glove	shoe	backpack			Shoe, Bag, and Accesso
sailboat unicycle motorcycle motorbike skateboard hot air balloon rocket bicycle Wehicle house building roof bridge church Building picture frame movie (disc) playing cards table cloth usical instrument guitar drum flute violin Musical Instrument telephone laptop computer remote mouse keyboard printer desktop copier radio Digital Product kite toy cars toy legos robot doll hair brush toner blush serum emulsion sunscreen Beauty and Skincare bottle spoon clock toothbrush basket pillow power outlet light switch Daily Necessity	sailboat unicycle motorcycle motorbike skateboard hot air balloon rocket bicycle Vehicle house building roof bridge church Building picture frame movie (disc) playing cards table cloth uusical instrument guitar drum flute violin Musical Instrument telephone laptop computer remote mouse keyboard printer desktop copier radio Digital Product kite toy cars toy legos robot doll hair brush toner blush serum emulsion sunscreen Beauty and Skinca bottle spoon clock toothbrush box chopping board ladder basket pillow power outlet light switch Daily Necessity	pen	pencil	fax machine	stapler				Stationery
unicycle motorcycle motorbike skateboard house building roof bridge church Building picture frame movie (disc) playing cards table cloth Artwork uusical instrument guitar drum flute violin Musical Instrument telephone laptop computer tablet ipad iphone cell phone Digital Product kite toy cars toy legos robot doll hair brush toner blush serum emulsion sunscreen Beauty and Skincare bottle plate clock toothbrush vase towel candle balloon Daily Necessity	unicycle motorcycle motorbike skateboard house building roof bridge church Building picture frame movie (disc) playing cards table cloth Artwork nusical instrument guitar drum flute violin Musical Instrument telephone laptop computer tablet ipad iphone cell phone Digital Product kite toy cars toy legos robot doll hair brush toner blush serum emulsion sunscreen Beauty and Skinca bottle plate cup bowl teapot fork knife spoon clock toothbrush vase towel candle balloon box chopping board ladder basket pillow power outlet light switch Daily Necessity	vehicle	car	van	truck	bus	train	boat	
house building roof bridge church Building picture frame movie (disc) playing cards table cloth Instrument guitar drum flute violin Musical Instrument telephone laptop computer keyboard printer desktop copier radio Digital Product kite toy cars toy legos robot doll hair brush toner blush serum emulsion sunscreen Beauty and Skincare bottle spoon clock toothbrush vase towel candle balloon box chopping board ladder basket pillow power outlet light switch Daily Necessity	house building roof bridge church Building picture frame movie (disc) playing cards table cloth nusical instrument guitar drum flute violin Musical Instrument telephone laptop computer tablet printer desktop copier radio Digital Product kite toy cars toy legos robot doll hair brush toner blush serum emulsion sunscreen Beauty and Skinca bottle spoon clock toothbrush vase towel candle balloon baily Necessity					hot air balloon	rocket	bicycle	Vehicle
picture frame movie (disc) playing cards table cloth Artwork usical instrument guitar drum flute violin Musical Instrument telephone laptop computer tablet ipad desktop copier radio Digital Product kite toy cars toy legos robot doll hair brush toner blush serum emulsion sunscreen Beauty and Skincare bottle spoon clock toothbrush vase towel candle balloon box chopping board ladder basket pillow power outlet light switch Daily Necessity	picture frame movie (disc) playing cards table cloth usical instrument guitar drum flute violin telephone laptop computer tablet ipad iphone cell phone remote mouse keyboard printer desktop copier radio kite toy cars toy legos robot doll hair brush toner blush serum emulsion sunscreen Beauty and Skinca bottle spoon clock toothbrush vase towel candle balloon box chopping board ladder basket pillow power outlet light switch Artwork Artwork Musical Instrument iphone cell phone cell phone radio Digital Product tablet ipad iphone cell phone radio Digital Product to doll Beauty and Skinca towel candle balloon balloon power outlet light switch	unicycle	motorcycle	motorbike	skateboard				
telephone laptop computer remote mouse keyboard printer desktop copier radio Digital Product kite toy cars toy legos robot doll hair brush toner blush serum emulsion sunscreen Beauty and Skincare bottle spoon clock toothbrush some box chopping board ladder basket pillow power outlet light switch Daily Necessity	nusical instrument guitar drum flute violin Musical Instrument telephone laptop computer remote mouse keyboard printer desktop copier radio Digital Product kite toy cars toy legos robot doll hair brush toner blush serum emulsion sunscreen Beauty and Skinca bottle spoon clock toothbrush vase towel candle balloon box chopping board ladder basket pillow power outlet light switch Musical Instrument Musical Instr	house	building	roof	bridge	church			Building
telephone laptop computer tablet ipad desktop copier radio Digital Product kite toy cars toy legos robot doll hair brush toner blush serum emulsion sunscreen Beauty and Skincare bottle spoon clock toothbrush vase towel candle balloon box chopping board ladder basket pillow power outlet light switch Daily Necessity	telephone laptop computer tablet jpad iphone cell phone mouse keyboard printer desktop copier radio Digital Product kite toy cars toy legos robot doll hair brush toner blush serum emulsion sunscreen Beauty and Skinca bottle plate cup bowl teapot fork knife spoon clock toothbrush vase towel candle balloon box chopping board ladder basket pillow power outlet light switch Daily Necessity	picture frame	movie (disc)	playing cards	table cloth				Artwork
remote mouse keyboard printer desktop copier radio Digital Product kite toy cars toy legos robot doll hair brush toner blush serum emulsion sunscreen Beauty and Skincare bottle spoon clock toothbrush box chopping board ladder basket pillow power outlet light switch Daily Necessity	remote mouse keyboard printer desktop copier radio Digital Product kite toy cars toy legos robot doll hair brush toner blush serum emulsion sunscreen Beauty and Skinca bottle spoon clock toothbrush vase towel candle balloon box chopping board ladder basket pillow power outlet light switch Daily Necessity	nusical instrument	guitar	drum	flute	violin			Musical Instrument
kite toy cars toy legos robot doll hair brush toner blush serum emulsion sunscreen Beauty and Skincare bottle plate cup bowl teapot fork knife spoon clock toothbrush vase towel candle balloon box chopping board ladder basket pillow power outlet light switch Daily Necessity	kite toy cars toy legos robot doll hair brush toner blush serum emulsion sunscreen Beauty and Skinca bottle plate cup bowl teapot fork knife spoon clock toothbrush vase towel candle balloon box chopping board ladder basket pillow power outlet light switch Daily Necessity								Di i ID I
hair brush toner blush serum emulsion sunscreen Beauty and Skincare bottle plate cup bowl teapot fork knife spoon clock toothbrush vase towel candle balloon box chopping board ladder basket pillow power outlet light switch Daily Necessity	hair brush toner blush serum emulsion sunscreen Beauty and Skinca bottle plate cup bowl teapot fork knife spoon clock toothbrush vase towel candle balloon box chopping board ladder basket pillow power outlet light switch Daily Necessity			•	<u> </u>	•	•	radio	Digital Product
bottle plate cup bowl teapot fork knife spoon clock toothbrush vase towel candle balloon box chopping board ladder basket pillow power outlet light switch Daily Necessity	bottle plate cup bowl teapot fork knife spoon clock toothbrush vase towel candle balloon box chopping board ladder basket pillow power outlet light switch		•	•					
spoon clock toothbrush vase towel candle balloon box chopping board ladder basket pillow power outlet light switch Daily Necessity	spoon clock toothbrush vase towel candle balloon box chopping board ladder basket pillow power outlet light switch Daily Necessity								Beauty and Skincare
box chopping board ladder basket pillow power outlet light switch Daily Necessity	box chopping board ladder basket pillow power outlet light switch Daily Necessity								
box cnopping board ladder basket pillow power outlet light switch	box chopping board ladder basket pillow power outlet light switch								Daily Necessity
	person Person		chopping board	ladder	basket	pillow	power outlet	ngnt switch	, ,

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-40:

"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.

577 The category is {}.

571

574

579

578 The corresponding instances are as follows:"

C Details of Prompt Generation

The specific instructions used in prompt generation are detailed in Figure 4. During the actual generation process, some of the prompts produced by GPT-40 did not meet the required criteria.

Table 5: Based on the categories, we employ GPT-40 to generate keyword associations and further enhanced the results by incorporating manually curated keywords.

Public Pacility fire extinguisher traffic sign street lamp street station Food and Beverage deible oil instant noodles pincapple apple don't milk orange approved chicken and the provided approved chicken and the provided approved chicken and the provided approved approved	The Third-level Category				Keywords			
Pool and Beverage Pool and Bev	Vehicle							oil tanker
Frod and Beverage plane either oil minant soedles principle apple of donut chicken monoiles and principle apple vegetable donut chicken monoiles harmbarge avocado sports drink power sports of the principle of t	Musical Instrument			digital piano harmonica				african drum erhu
Food and Beverage Proposed and Beverage Redical Supply Redi	Public Facility	fire extinguisher	traffic sign	street lamp	street	station		
Book Yearhook Section Sectio	Food and Beverage	pineapple apple	milk donut	orange durian	avocado sports drink	can canned health products	juice egg	
Book notebook magazine dictionary	Medical Supply					first aid kit	medication	medicine bottle
Furniture	Book					encyclopedia	atlas	pamphlet
Home Appliance enicrophone television oven eligible oven beliably of the state of liberal to the state	Furniture	barber chair	office chair	bathroom mirror	chair	sofa	bean bag chair dining table	children's chair bed
Building house apartment building duplex house church statute of liberty effect lower eiffel tower eiffel tower of piss pyramid statue of liberty eiffel tower ei	Home Appliance	microphone	refrigerator	hair dryer	humidifier			robot vacuum curling iron
Digital Product Smart robot printer cancorder vintage camera camera monitor desktop computer smartwarch fines tracker monitor monito	Amphibian					Surinam toad	alpine newt	glass frog
Digital Product Smartwatch vintage camera monitor drone palptop projector fitness tracker	Building						castle	golden gate brid
Stationery glue stick stapler crayon ballpoint pen floppy disk eraser Daily Necessity birdcage glass jar vase hanger electric saw mop broon comb Plant cactus coconut tree mint rose sunflower tulip catus lavender choker gold barry bowl toothbrush shower gel clock kitchen knil electric saw mop broon comb Plant mint rose sunflower tulip cactus lavender choker gold barry bowl toothbrush shower gel clock kitchen knil electric saw mop broon comb Plant mint rose sunflower tulip cactus lavender hair accessory gold bar necklace sunflower tulip cactus lavender choker gold bar necklace and the pendant brooch anklet locket Beauty and Skincare perfume blush cye shadow facial serum erase sunflower eramic craft carmic craft wood carving classical bust serum mascara lipstick. Artwork bouquet of flowers clay sculpture ceramic craft down jacket coat skirt shorts vest shirt down jacket coat skirt shorts vest shorts west shorts west barbell dumbell Clothing pants shirt down jacket coat skirt shorts vest shorts west short probt motorbik toy majec cube poop emoil ite handbag sandals shoes luggage purse face poop emoil short bush barbapa do geases sandals backpack cap life down jacket coat skirt shorts vest short probt motorbik toy majec cube poop emoil ite handbag sandals shoes luggage purse face of the publish down jacket carmas shoes shorts shoes sandals backpack cap life down jacket carmas shoes luggage purse face of the publish down jacket carmas shoes luggage purse face of the publish down jacket carmas shoes luggage purse face of the publish down jacket carmas shoes luggage purse face of the publish down jacket carmas shoes luggage purse face of the publish down jacket carmas shoes luggage purse face of the publish down jacket carmas shoes luggage purse face of the publish down jacket carmas shoes luggage purse face of the publish down jacket carmas shoes luggage purse face of the publish down jacket carmas shoes shoes luggage purse face of the publish down jacket carmas shoes shoes shoes shorts sheet down jacked turtle	Digital Product	printer	camcorder	camera	smart camera	laptop	mobile phone	tablet walkie-talkie
Daily Necessity Daily Necessity Daily Neces	Insect	shrimp	crab	ant	grasshopper	butterfly		
Daily Necessity Daily Necessity glass jar electric saw mop broom bowl frying pan baby bottle kitchen knit	Stationery					tape measure	scissors	compass
Jewelry	Daily Necessity	birdcage glass jar	alarm clock vase	spoon hanger	bowl soap dish	toothbrush	shower gel	
Beauty and Skincare	Plant							maple leaf
Beauty and Skincare blush eye shadow facial serum emulsion serum mascara lipstick	Jewelry	tiara	crown	stud	chain	gemstone	choker	beaded bracele hairpin
Clothing dress pants baby clothes clothing jeans sweatshirt T-shirt socks pants babry clothes down jacket coat skirt shorts vest tennis adjustable bench treadmill skateboard skateboard backpack soccer sleeping bag baseball flamingo flo dumbbell skateboard skateboard backpack soccer sleeping bag baseball flamingo flo dumbbell short treadmill skateboard skateboard backpack soccer sleeping bag baseball flamingo flo dumbbell short treadmill skateboard skateboard backpack soccer sleeping bag baseball flamingo flo dumbbell short treadmill skateboard skateboard shores luggage purse fancy boot belt sneaker sandals shoes luggage purse fancy boot hard backpack cap tie handbag sandals sandals shoes luggage purse fancy boot motorbike toy magic cube poop emoji sloth plushie bear plushie bear plushie red cartoor feeve figuri doll plushie short puppy monkey kitten dolphin french bulldog french	Beauty and Skincare							
Sports Equipment dennis adjustable bench treadmill shorts ball tent treadmill backpack barbell tent treadmill backpack barbell tent treadmill backpack barbell backpack backpack barbell shorts backpack barbell shorts backpack barbell shorts backpack cap tie backpack backpack backpack cap backpack backpack cap backpack backpack cap backpack backpack cap backpack cap backpack backpack cap backpack backpack cap backpack backpack backpack cap backpack backpack backpack cap backpack b	Artwork					stone carving	catstatue	mugskulls
Sports Equipment adjustable bench treadmill skateboard barbell dumbbell dumbbell skateboard barbell dumbbell statebard dumbbell dumbbell statebard barbell dumbbell statebard dumbbell skateboard barbell dumbbell statebard d	Clothing							
hoe, Bag and Accessory glasses sandals backpack cap tie fancy boot handbag sandals sandals fancy boot handbag sandals sandals fancy boot fancy boot fancy boot fancy boot fancy foot fancy f	Sports Equipment	adjustable bench	knee pad	backpack	soccer			badminton flamingo float
Toy robot minion smart robot robot toy toy wolf plushie bear plushie red cartoor wolf plushie wolf plushie doll Eeve figuri Mammal rabbit fox wolf Siamese cat together dog raccoon lion panda alpaca puppy monkey kitten dolphin French bulldog Reptile cobra gecko rattlesnake crocodile snake lizard utrute sea turtle soft-shelled turtle snake lizard Bird heron pigeon toucan parrot stork flamingo penguin woodpecker nightingale duck turkey chicken crow eagle peacock swallow owl kingfisher peacock bird canary sparrow rooster Fish shark tropical fish jellyfish goldfish perch eel monkfish	hoe, Bag and Accessory	glasses	sandals	shoes	luggage purse	fancy boot	belt	
Mammal panda alpaca elephant puppy Ilama monkey tiger kitten dog holphin raccoon French bulldog lion French bulldog Reptile cobra turtle gecko rattlesnake crocodile chameleon lizard alligator iguana Bird heron woodpecker nightingale peacock toucan owl kingfisher peacock parrot owl kingfisher hawk dove anchovy Fish shark tropical fish jellyfish goldfish perch eel monkfish	Toy	robot	motorbike toy	magic cube	poop emoji	sloth plushie	bear plushie	balloon red cartoon Eevee figurine
Bird heron pigeon toucan parrot stork flamingo penguin woodpecker nightingale peacock swallow bird canary sparrow rooster Fish shark tropical fish jellyfish goldfish perch lizard soake lizard stork flamingo penguin chicken crow eagle duck turkey chicken crow eagle anchovy sparrow rooster	Mammal	panda	elephant	llama	tiger	dog	raccoon	
Bird woodpecker nightingale duck turkey chicken crow eagle peacock swallow owl kingfisher hawk dove anchovy bird canary sparrow rooster Fish shark tropical fish jellyfish goldfish perch eel monkfish	Reptile						alligator	iguana
	Bird	woodpecker peacock	nightingale swallow	duck owl	turkey kingfisher	chicken	crow	eagle
	Fish			jellyfish				

Therefore, we instructed GPT-40 to generate multiple prompts for each image, and then manually selected those that best matched the intended scenarios. Figure 13 presents the results generated by different methods in this study, along with their corresponding prompts.

85 D Additional Discussions and Details of Model Performance

D.1 Additional Discussions

Analysis of The First-Level Category The primary categories are divided into photorealistic and non-photorealistic. Table 6 and Figure 8 present the performance of different methods on these two categories across three evaluation dimensions. The results show that: (1) Subject Preservation: Almost all methods perform better on photorealistic categories than on non-photorealistic ones. We speculate that this is because, when referencing subjects from non-photorealistic categories, these methods tend to generate photorealistic images based on the prompt, which results in lower subject consistency. (2) Prompt Following: The performance gap between photorealistic and non-photorealistic categories is relatively small. This can be attributed to the fact that CLIP-T focuses primarily on the semantic information of the image. As long as the generated subject matches the category and general appearance described in the prompt, the CLIP-T score will not be significantly reduced. (3) Image quality: There is little difference in performance between photorealistic and non-photorealistic categories. This indicates that the distinction between these two categories does not affect the quality of image generation, and the HPSv2 metric does not show a preference for either category.

Analysis of The Second-Level Category The secondary categories under both the realistic and non-realistic primary categories are further subdivided into objects, humans, and animals. Table 7 and Figure 9 present the performance of various methods across these three dimensions for both realistic and non-realistic categories. The results demonstrate that, irrespective of whether the primary category is realistic or non-realistic, the scores for the subject preservation dimension are consistently lower for the human category across nearly all models. As detailed in Table 8, this phenomenon can be attributed to the distribution of difficulty levels within the human category, where the proportions of simple, medium, and hard cases are 1.96%, 50.98%, and 47.06%, respectively. In contrast, the object and animal categories exhibit a higher proportion of subjects at the simple difficulty level and a lower proportion at the hard difficulty level, which likely contributes to their relatively higher subject preservation scores.

Implications for Technical Approaches (1) Figure 10 shows that, as base models and model architectures are updated, the performance boundary of these models consistently expands outward. Table 9 presents all the base models used by each method. It can be observed that the top-performing methods consistently employ relatively recent text-to-image base models. For instance, UNO utilizes FLUX as its foundational model. This observation suggests that the adoption of advanced text-toimage base models is a critical factor in enhancing performance on subject-driven T2I tasks. (2) Historically, fine-tuning methods have generally outperformed encoder-based approaches in terms of subject preservation. This advantage is attributed to their ability to better retain the original text-image conditional distribution by fine-tuning on images of the specified subject. In contrast, encoder-based methods often encounter interference during feature injection, which can hinder precise prompt alignment. However, with the development of more advanced encoding techniques, the adoption of larger and more powerful base models, and the availability of extensive training datasets, encoder-based methods have demonstrated significantly improved performance. From an application standpoint, fine-tuning methods require substantial computational resources for optimization and often exhibit limited generalization capabilities. In contrast, encoder-based methods are less constrained by these limitations, making them more practical for future applications. Nevertheless, our analysis indicates that current encoder-based methods still face challenges in accurately reconstructing subjects with high-frequency details in images. This limitation may stem from the characteristics of commonly used image encoders, such as CLIP, which tend to prioritize semantic information over fine-grained details. Consequently, future research should focus on enhancing the restoration of challenging subject details.

D.2 Details of Model Performance

In this section, we present the detailed evaluation results for each metric across all models. To comprehensively evaluate the effectiveness of different metrics for assessing subject consistency, we calculated multiple metrics for each method. The detailed results are presented in Table 10, 11, 12. In section 5, we present the performance of all methods across images with different difficulty

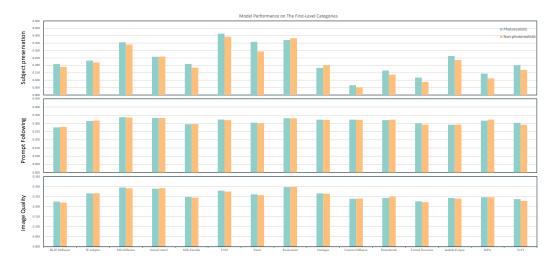


Figure 8: Comparison of bar charts for DSH-Bench scores in different first-level categories.

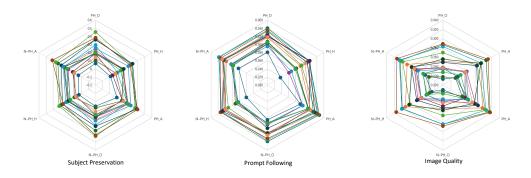


Figure 9: Comparison of radar charts for DSH-Bench scores in different second-level categories.

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 I	Preservation [†]	Prompt	Following [†]	Image (Quality [†]
	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	0.338	0.336	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	0.414	0.394	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	0.297	0.298
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
λ -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
Aver.	0.236	0.217	0.313	0.312	0.257	0.256

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			Subje	ct Preservation		
Wellou	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	0.453	0.312	0.361	0.428	0.315	0.365
Emu2	0.358	0.326	0.387	0.305	0.266	0.285
RealCustom++	0.383	0.291	0.396	0.415	0.26	0.412
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
λ -Eclipse	0.252	0.266	0.101	0.236	0.221	0.256
HiPer	0.143	0.283	0.195	0.126	0.079	0.098
NeTI	0.143	0.083	0.193	0.120	0.079	0.201
Aver.	0.193	0.139	0.262	0.104	0.130	0.223
Tiver.	0.211	0.100		npt Following	0.100	0.223
Method	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	0.201	0.237	0.322	0.317	0.319	0.232
MS-Diffusion	0.317	0.234	0.322	0.338	0.332	0.317
OminiControl	0.340	0.319	0.344	0.334	0.332	0.337
SSR-Encoder		0.319	0.344	0.334	0.33	
	0.302					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	0.333	0.333
OmniGen	0.320	0.320	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
λ -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
Aver.	0.315	0.294	0.322	0.313	0.307	0.316
Method			Im	age 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++	0.289	0.312	0.317	0.288	0.316	0.314
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
λ -Eclipse	0.234	0.257	0.263	0.23	0.254	0.262
HiPer	0.237	0.256	0.203	0.238	0.255	0.202
NeTI	0.228	0.244	0.261	0.218	0.24	0.249

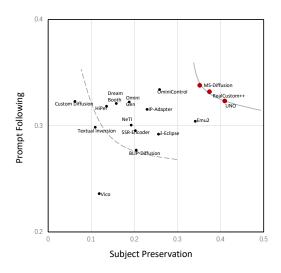


Figure 10: Pareto front diagram illustrating model performance across both subject and prompt dimensions. The red points in the diagram represent the current Pareto-optimal solutions.

Table 8: Subject hard level distribution under the second category

	Photorealistic											
Benchmark		Object			Human			Animal				
DreamBench++ DSH-Bench Benchmark DreamBench DreamBench++	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			
	Non-photorealistic											
Benchmark		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			

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.

40 E Implementation Details

641 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

646

Evaluation Instruction Figure 11 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 11.

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	0.298	0.110
UNO	FLUX.1-dev	0.409	0.323	0.278	0.109
MS-Diffusion	SDXL	0.352	0.338	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
λ -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		Subjec	t Pres	ervation		Prompt	Following	Ima	ge Quality	
	C-B-I↑	C-L-I↑	D-I↑	D-v2-I↑	SICS↑	С-В-Т↑	C-L-T↑	ImageReward↑	PickScore↑	HPSv2↑
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	0.820	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	0.336	0.284	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	0.716	0.409	0.317	0.259	0.725	0.602	0.304
Emu2	0.838	0.818	0.737	0.704	0.360	0.291	0.235	0.463	0.599	0.272
RealCustom++	0.794	0.770	0.746	0.698	0.377	0.325	0.278	0.813	0.601	0.316

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		Subjec	t Pres	ervation		Prompt	Following	Ima		
	C-B-I↑	C-L-I↑	D-I↑	D-v2-I↑	SICS↑	С-В-Т↑	C-L-T↑	ImageReward↑	PickScore↑	HPSv2↑
BLIP-Diffusion	0.836	0.809	0.691	0.664	0.216	0.279	0.225	0.260	0.591	0.249
IP-Adapter	0.846	0.845	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	0.339	0.285	0.729	0.593	0.309
OminiControl	0.761	0.780	0.551	0.566	0.268	0.336	0.284	0.793	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	0.410	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	0.714	0.706	0.380	0.330	0.280	0.710	0.594	0.314

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		Subjec	t Prese	ervation		Prompt	Following	Ima	ge Quality	
	C-B-I↑	C-L-I↑	D-I↑	D-v2-I↑	SICS↑	С-В-Т↑	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	0.338	0.287	0.705	0.595	0.294
OminiControl	0.721	0.736	0.462	0.461	0.258	0.334	0.288	0.787	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	0.409	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	0.298
OmniGen	0.696	0.678	0.436	0.326	0.188	0.322	0.274	0.586	0.592	0.265
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
λ -Eclipse	0.852	0.833	0.676	0.638	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														Subjec	Prese	rvation													
c.iiou	VE	MI	PUF	FB	MS	во	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	0.268	0.192	0.205	0.181	0.202	0.202	0.194
IP-Adapter			0.183	0.217		0.137		0.228	0.208	0.300	0.216	0.123	0.257			0.189					0.262	0.218		0.190		0.225	0.195	0.221	
MS-Diffusion		0.368			0.378			0.361	0.592		0.346	0.244	0.353				0.307		0.449		0.433	0.336				0.337	0.337	0.302	
OminiControl SSR-Encoder	0.187	0.329		0.312	0.314	0.172		0.335	0.275		0.296		0.296				0.296	0.293	0.258		0.336	0.324		0.208		0.208	0.135	0.082	
UNO		0.411		0.191	0.461	0.130			0.192		0.163	0.146	0.138				0.138		0.204					0.197		0.213	0.159	0.228	
Emu2		0.281			0.382	0.302												0.348				0.300		0.310		0.510	0.308	0.311	
RealCustom++		0.421	0.371	0.340	0.328	0.228		0.395	0.500	0.406		0.327			0.432				0.389	0.488	0.369	0.417				0.390	0.306	0.270	
OmniGen		0.123		0.249	0.161			0.133	0.058		0.192	0.096							0.202		0.203	0.180				0.144	0.209	0.215	
Custom Diffusion		0.083	0.060	0.086	0.058	0.037		0.051	0.083	0.181		0.063	0.063				0.082					0.055				0.060	0.014	0.015	
DreamBooth		0.180	0.206	0.206	0.151	0.077		0.232	0.200	0.242		0.190							0.131		0.167	0.186			0.190		0.048	0.035	
		0.081		0.116		0.022					0.059	0.117						0.101				0.134			0.219			0.078	
λ-Eclipse HiPer	0.280	0.215		0.215	0.214	0.170	0.346	0.248	0.283	0.312			0.211	0.249					0.298		0.309	0.238		0.227	0.272	0.231	0.261 0.092	0.273	
NeTI										0.201																	0.092		
	0.199	0.154	0.140	0.193	0.179	0.123	0.233	0.200	0.292	0.275	0.144	0.190	0.102		_		0.132	0.100	0.143	0.193	0.193	0.220	0.202	0.232	0.202	0.233	0.132	0.139	0.210
Method	VE	MI	PUF	FB	MS	ВО	FUR	HA	AM	BU	DP	IN	ST	DN	pt Follo	JE	BS	AR	CL	SE	SBA	то	MA	RE	BI	FI	HF	FA	AC
BLIP-Diffusion	0.271	0.283	_	0.282	0.275	0.244	0.269	0.286	0.307	0.285	0.273		0.283	0.286	0.294	_		0.272	0.281	0.296	0.281	0.294	0.291	0.287		0.287	0.241	0.223	0.245
IP-Adapter		0.265		0.282	0.273	0.244	0.209	0.200	0.307	0.283	0.273	0.299	0.320			0.200	0.307		0.281		0.321	0.294	0.291	0.303		0.287	0.241	0.223	0.243
MS-Diffusion				0.339										0.344								0.353			0.345				
OminiControl				0.337				0.336					0.328		0.339										0.344			0.321	
SSR-Encoder	0.290			0.301	0.295				0.303	0.294	0.297		0.306		0.311	0.300	0.307		0.304		0.299	0.309		0.281		0.293	0.267	0.249	
UNO				0.328	0.325			0.331									0.320				0.334					0.321	0.295	0.299	
Emu2		0.316	0.307	0.303	0.309	0.285	0.308	0.310					0.310		0.319		0.296		0.305		0.308	0.315				0.309	0.292	0.276	
RealCustom++		0.339		0.336	0.333				0.352				0.323				0.323			0.338		0.344				0.325	0.312		
OmniGen Custom Diffusion		0.318	0.312	0.333	0.312				0.337							0.310	0.323	0.317			0.327	0.328		0.311	0.331	0.317		0.317	
DreamBooth			0.315		0.321	0.313												0.320			0.325				0.330		0.312		
Textual Inversion		0.307		0.308	0.304	0.275	0.293		0.323				0.300				0.302		0.299					0.301	0.316		0.276	0.290	
λ-Eclipse	0.280	0.302	0.286	0.296	0.298	0.278			0.312		0.290		0.303	0.310					0.292						0.297		0.272	0.274	0.248
HiPer				0.316						0.312										0.331					0.320		0.305		
NeTI	0.308	0.314	0.311	0.311	0.314	0.276	0.300	0.312	0.302	0.307	0.302	0.310	0.300	0.302	0.312	0.295	0.307	0.299	0.307	0.315	0.315	0.312	0.320	0.309	0.311	0.300	0.281	0.282	0.266
Method														Ima	ge Qua	lity													
	VE	MI	PUF	FB	MS	во	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	0.253	0.203	0.213	0.222	0.175	0.168		0.190	0.274		0.196	0.256	0.195			0.206	0.205		0.216		0.227	0.243	0.263	0.246		0.263	0.234	0.229	0.233
IP-Adapter		0.248	0.241		0.218							0.277						0.279	0.255			0.284		0.266		0.294	0.290	0.293	
MS-Diffusion		0.279		0.290	0.260	0.269	0.282	0.279	0.310	0.306	0.279	0.305	0.271	0.287			0.275		0.288					0.288		0.309	0.303	0.311	
OminiControl SSR-Encoder		0.277	0.280	0.283	0.264	0.286		0.274	0.313	0.299	0.280	0.287	0.273	0.282		0.277	0.274		0.284	0.277	0.290	0.303	0.314	0.283		0.292	0.292	0.298	0.300
UNO UNO	0.275	0.226	0.239	0.239	0.202	0.204	0.216	0.225	0.290	0.264	0.223	0.289	0.221	0.230					0.243		0.242	0.259		0.284		0.282	0.281	0.254	0.287
Emu2			0.277	0.243	0.248	0.214	0.232	0.234	0.288	0.276	0.234		0.246	0.248		0.234	0.202		0.272		0.279	0.293		0.271		0.290	0.282	0.287	0.287
RealCustom++		0.279		0.291					0.307			0.313						0.297		0.290				0.301	0.319		0.311		
OmniGen	0.278	0.239	0.253	0.263	0.237	0.240	0.251	0.242	0.265	0.286	0.253	0.255	0.251	0.251	0.271	0.243	0.258	0.271	0.241	0.256	0.260	0.266	0.286	0.249	0.281	0.263	0.300	0.292	0.281
Custom Diffusion	0.255	0.234	0.236	0.235	0.223	0.221		0.230	0.245	0.249	0.226	0.239	0.219			0.228	0.231		0.237		0.241	0.249		0.247		0.247	0.235	0.237	0.238
DreamBooth	0.266			0.239					0.222		0.228	0.250						0.252		0.237					0.252			0.255	
Textual Inversion	0.253	0.214	0.219	0.226	0.212	0.178	0.208	0.214	0.261	0.248	0.205	0.233	0.206	0.205		0.207	0.211		0.218	0.217	0.228	0.231	0.252			0.236	0.227	0.232	
λ-Eclipse	0.276	0.220		0.235	0.206	0.222		0.222	0.277	0.242	0.234	0.248	0.229	0.236			0.227		0.238		0.237	0.247	0.266	0.245		0.267	0.263	0.256	0.245
HiPer NaTI						0.211			0.287									0.239							0.271		0.257	0.255	
NeTI	0.262	0.224	0.234	0.233	0.213	0.186	0.215	0.218	0.249	0.248	0.211	0.246	0.214	0.221	0.247	0.214	0.224	0.246	0.228	0.234	0.241	0.245	0.266	0.250	0.257	0.260	0.243	0.240	0.250

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														Subjec	t Prese	rvation													
c.mou	VE	MI	PUF	FB	MS	во	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			0.208	0.146	0.153	0.210	0.180		0.185	0.139		0.108	0.208			0.181			0.375	0.221	0.177	0.219		0.192		0.180	0.221	0.194
IP-Adapter				0.283				0.210										0.175						0.233				0.200	
MS-Diffusion		0.410		0.375	0.396		0.363			0.292			0.233			0.296	0.310		0.438	0.742	0.426						0.302		
OminiControl	0.290	0.370	0.367	0.333	0.377	0.275	0.181		0.292	0.265	0.269		0.237				0.362			0.633	0.342		0.273		0.190	0.283	0.158	0.092	
SSR-Encoder		0.168					0.179											0.175						0.156			0.145		
UNO		0.448	0.537			0.428	0.415			0.377					0.317		0.481			0.675							0.315		0.338
Emu2	0.350	0.310	0.337	0.315		0.303	0.303		0.192	0.267	0.186		0.379					0.158	0.305	0.433		0.187	0.294			0.375	0.251	0.250	
RealCustom++ OmniGen		0.425		0.431	0.508	0.472	0.436			0.300										0.450							0.265		
Custom Diffusion		0.130	0.108	0.246	0.073	0.094					0.094		0.050				0.236			0.373		0.002	0.203				0.178	0.000	
DreamBooth		0.072					0.017			0.134			0.050														0.023		
		0.090	0.179	0.215	0.101	0.067	0.111			0.164	0.033		0.103				0.134	0.058	0.058	0.368							0.003		
λ-Eclipse			0.225			0.261	0.194		0.250				0.003			0.192				0.517			0.270		0.232				
HiPer								0.114		0.162								0.200						0.031					
NeTI																		0.317											
	0.101	0.200	0.207	0.140	0.154	0.104	0.10)	0.150	0.275	0.101	0.142	0.154	0.117		pt Folle		0.140	0.517	0.075	0.172	0.100	0.157	0.210	0.251	0.102	0.150	0.150	0.175	0.151
Method	VE	MI	PUF	FB	MS	ВО	FUR	HA	AM	BU	DP	IN	ST	DN	PL	JE	BS	AR	CL	SE	SBA	ТО	MA	RE	BI	FI	HF	FA	AC
BLIP-Diffusion	0.290	0.287	0.291	0.293	0.291	0.279	0.275		0.281	0.280	0.292		0.278			0.282	0.287			0.285	0.292		0.283		0.274				
IP-Adapter		0.336	0.316		0.335	0.304	0.307		0.326	0.304	0.313		0.321			0.320				0.319			0.321		0.311				
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
OminiControl	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	0.293	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	0.296	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			0.309	0.304	0.276	0.308	0.323	0.312	0.273	0.297	0.312	0.304	0.297	0.292	0.243	
RealCustom++	0.319		0.326		0.340	0.328	0.331		0.317	0.326	0.309	0.332				0.325	0.329			0.336	0.338		0.338		0.339		0.330	0.321	0.338
OmniGen						0.318		0.305		0.317		0.322					0.338			0.336				0.327			0.326		
Custom Diffusion	0.315			0.326	0.326	0.313	0.324		0.322	0.323	0.308		0.311		0.314		0.323			0.329	0.321		0.331		0.325		0.316		
DreamBooth		0.337		0.331	0.330		0.322	0.317		0.322	0.310	0.317					0.316						0.334		0.323		0.324		0.324
	0.290		0.287	0.293		0.287	0.291	0.293		0.301								0.266		0.314			0.299		0.298		0.290		
λ-Eclipse	0.280	0.313		0.300	0.293	0.290	0.295		0.301	0.278	0.281		0.302	0.310			0.306			0.335			0.305		0.306		0.286		
HiPer NeTI		0.339	0.317			0.317			0.327									0.220					0.336		0.320		0.319		
Nell	0.297	0.327	0.303	0.300	0.317	0.292	0.274	0.291	0.302	0.293	0.272	0.300	0.301				0.300	0.230	0.291	0.323	0.200	0.270	0.299	0.290	0.293	0.301	0.263	0.227	0.290
Method															ge Qua														
	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		0.220	0.221	0.189	0.224	0.220		0.212	0.184	0.225		0.246		0.229		0.224		
IP-Adapter		0.255	0.239	0.263	0.230	0.231	0.227		0.309	0.279	0.245		0.259	0.258	0.220	0.278	0.264			0.229	0.251		0.294				0.283	0.293	
MS-Diffusion		0.287	0.268	0.296	0.269	0.273	0.269		0.308	0.304	0.267		0.288	0.294			0.298		0.284	0.283	0.289		0.307		0.296		0.293	0.307	
OminiControl		0.290	0.265	0.287	0.269		0.278			0.298	0.277		0.284				0.285		0.282	0.281		0.270					0.297		
SSR-Encoder	0.258	0.237	0.216	0.241	0.204	0.206	0.213		0.288	0.260		0.242	0.232		0.217	0.245	0.237	0.227	0.229	0.210	0.241		0.271	0.287	0.266		0.260	0.248	
UNO Emu2											0.243				0.223					0.228				0.315		0.311			
RealCustom++		0.201		0.259	0.269	0.239	0.227			0.269			0.246				0.252			0.233			0.269				0.268		
OmniGen	0.304		0.283			0.262							0.248				0.282			0.268			0.281				0.303		
Custom Diffusion						0.202	0.241			0.257										0.245				0.245			0.276		
DreamBooth		0.249	0.246			0.221	0.241		0.240	0.237	0.221		0.213						0.233	0.245	0.233		0.269		0.262		0.249		
Textual Inversion	0.227		0.240			0.198	0.233		0.240	0.270			0.229				0.230		0.240	0.243	0.243		0.230		0.202			0.249	
λ-Eclipse	0.250						0.212		0.264			0.252			0.212			0.208					0.263		0.249		0.252		
HiPer						0.226		0.227		0.245				0.238				0.193					0.203		0.265		0.252		
NeTI		0.244																0.208								0.277			

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.

	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	0.433	0.414	0.379	0.359	0.418	0.349								
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								
λ -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								
Aver.	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	0.339	0.341	0.324	0.338	0.341
OminiControl	0.338	0.334	0.337	0.329	0.326	0.342
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++	0.343	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
λ -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
Aver.	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	0.293	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++	0.300	0.298	0.295	0.284	0.301	0.307
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
λ -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
Aver.	0.260	0.255	0.251	0.258	0.260	0.257

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.

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 SICS, CLIP-T, and HPSv2, respectively.

Method	Sub	ject Preserva	tion	Pr	ompt Follow	ing	1	mage Qualit	y
	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	0.340	0.339	0.335	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	0.469	0.405	0.383	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	0.281	0.300	0.303
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
λ -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
Aver.	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
λ -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

MLLM Projects (1) Qwen2.5-VL-7B: https://huggingface.co/Qwen/Qwen2.5-VL-7B-Instruct (2) Implementation Framework: https://github.com/hiyouga/LLaMA-Factory

Why we use *Kendall's* τ *value* and *Spearman correlation coefficient value*: The scenario is as follows: We have multiple metrics scoring the same dataset. These metrics may have different value ranges. The ground truth scores are provided by human annotators. We want to measure the correlation between each metric and the human scores. The following are some commonly used correlation evaluation metrics:

1. Pearson Correlation Coefficient

- Advantages: Measures linear correlation between two continuous variables; simple to compute.
- **Disadvantages:** Only suitable for linear relationships;sensitive to outliers; requires interval or ratio data; variables should be on the same scale.
- **Applicability:** Use if both our metric and human scores are continuous and linearly related. If scales differ, standardize first.

2. Spearman Rank Correlation Coefficient

- Advantages: Measures monotonic relationships; does not require linearity; robust to outliers; works with different scales and ordinal data.
- **Disadvantages:** Only captures monotonic relationships; some information loss due to ranking.

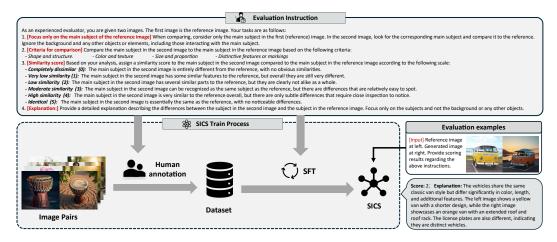


Figure 11: **The training process of SICS.** We constructed and annotated a dataset specifically tailored for subject consistency determination, and subsequently trained models using this dataset.

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

3. Kendall Rank Correlation Coefficient

- 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

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690 691

693

698

699

700

701

702

703

704

- 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.
- 692 Consequently, we choose Kendall's \u03c4 value and Spearman correlation coefficient value.

F More Generation Examples

- Figure 12 shows the generation examples of different methods across different difficulty levels.
- Figure 13 shows the generation examples of different methods across different prompt scenarios.
- The blue block highlights encoder-based methods, and the green block highlights fine-tuning-based methods.

G Limitations

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

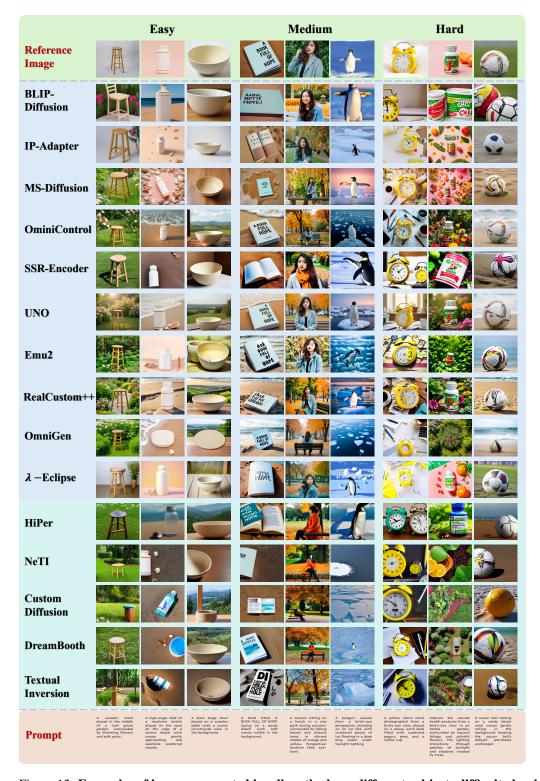


Figure 12: Examples of images generated by all methods on different subjects difficulty level.

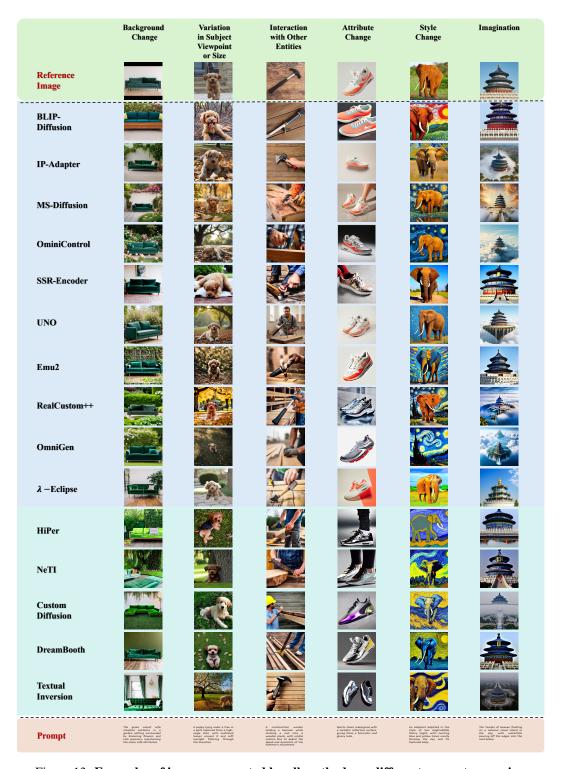


Figure 13: Examples of images generated by all methods on different prompt scenarios.

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: See Section 3.1, 4.2, 5.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the
 contributions made in the paper and important assumptions and limitations. A No or
 NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals
 are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: We discuss in Section G.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was
 only tested on a few datasets or with a few runs. In general, empirical results often
 depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: We do not have theoretical result.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and crossreferenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: See Section 4.1,E, 3.1.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
- (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

811 Answer: Yes

813

815

816

817

818

821

822

823

824

825

826

827

828

829

830

831

832

833

834 835

836 837

838

840

841

842

843 844

845

846

847

848 849

850

851

852

853

854

855

856

859

860

861

Justification: We provide the code and model in supplementay material. We are open-sourcing full data.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how
 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new
 proposed method and baselines. If only a subset of experiments are reproducible, they
 should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: See Section 4.1,E.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [No]

Justification: Most of our experiments are evaluating existing Diffusion models.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).

- It should be clear whether the error bar is the standard deviation or the standard error
 of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

862

863

864

865

866

867

868

869

870

871

872

873 874

875

876

877

878

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

Justification: See Section E.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: We ensured that each image's copyright status was verified for academic suitability, and manual inspection was conducted to confirm that the prompts used were free of defects.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a
 deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader Impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: See Section G.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.

- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [Yes]

Justification: We manually reviewed all the data in the benchmark to avoid unsafe images and prompts.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We ensured that the sources of both each image and the corresponding model meet the required standards.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the
 package should be provided. For popular datasets, paperswithcode.com/datasets
 has curated licenses for some datasets. Their licensing guide can help determine the
 license of a dataset.

- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New Assets

965

966

967

968

969

970

971

972

973

974

975

976

977

978

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

995

996

997

998

999

1001

1002

1003

1004

1005

1006

1007

1008

1009

1010

1011

1012

1013

1014

1015

1016

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Justification: We curated new annotations.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and Research with Human Subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: This paper does not involve crowdsourcing nor research with human subjects. Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: This paper does not involve crowdsourcing nor research with human subjects. Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.