

---

# Why Settle for One? Text-to-ImageSet Generation and Evaluation

---

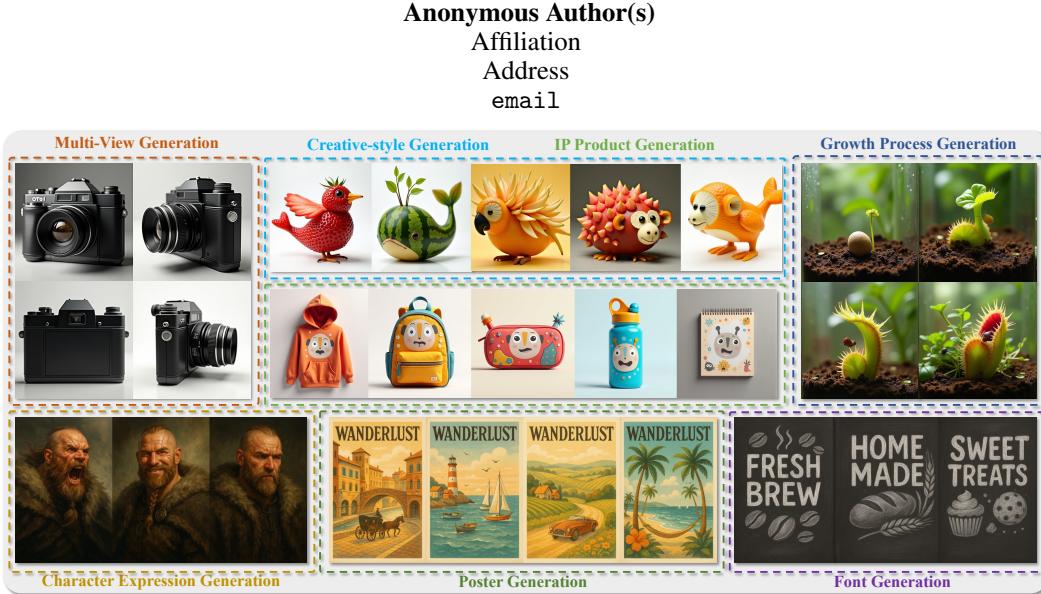


Figure 1: Illustration of diverse applications of Text-to-ImageSet generation across multiple domains. All examples are from the proposed benchmark. Images in the upper are generated by our AutoT2IS, while images in the lower are generated by AutoT2IS combined with the commercial model GPT-4o.

## Abstract

1 Despite remarkable progress in Text-to-Image models, many real-world applica-  
2 tions require generating coherent image sets with diverse consistency requirements.  
3 Existing consistent methods often focus on a specific domain with specific aspects  
4 of consistency, which significantly constrains their generalizability to broader appli-  
5 cations. In this paper, we propose a more challenging problem, Text-to-ImageSet  
6 (T2IS) generation, which aims to generate sets of images that meet various con-  
7 sistency requirements based on user instructions. To systematically study this  
8 problem, we first introduce **T2IS-Bench** with 596 diverse instructions across 26  
9 subcategories, providing comprehensive coverage for T2IS generation. Building  
10 on this, we propose **T2IS-Eval**, an evaluation framework that transforms user  
11 instructions into multifaceted assessment criteria and employs effective evaluators  
12 to adaptively assess consistency fulfillment between criteria and generated sets.  
13 Subsequently, we propose **AutoT2IS**, a training-free framework that maximally  
14 leverages pretrained Diffusion Transformers’ in-context capabilities to harmonize  
15 visual elements to satisfy both image-level prompt alignment and set-level visual  
16 consistency. Extensive experiments on T2IS-Bench reveal that diverse consistency  
17 challenges all existing methods, while our AutoT2IS significantly outperforms cur-  
18 rent generalized and even specialized approaches. Our method also demonstrates  
19 the ability to enable numerous underexplored real-world applications, confirming  
20 its substantial practical value. All our data and code will be publicly available.

21    **1 Introduction**

22    Recent advancements in Text-to-Image (T2I) models [3, 4, 11, 24, 31, 37, 39] have substantially  
23    enhanced the ability to generate visually compelling and semantically faithful images. However, as  
24    illustrated in Figure 1, most practical scenarios, *e.g.*, product design, process illustrations, or character  
25    creation, often require not a single image, but a coherent image set. Such image set must maintain  
26    ***set-level visual consistency*** with varying degrees of ***identity preservation*** [19], ***uniform style*** [16],  
27    and ***logical coherence*** [41]. This novel requirement poses huge challenges for current T2I models,  
28    which are primarily focused on image-level prompt alignment rather than set-level visual consistency.

29    Previous attempts addressing this coherence have largely concentrated on specific facets of consistency [25, 41, 48] or specific combinations of consistency types [5, 30]. For instance, consistent  
30    character generation approaches [25, 29, 45] aim to preserve character identity across various contexts,  
31    while methods such as StoryDiffusion [51] are tailored for generating stylistically coherent comic  
32    sequences with localized identity preservation. Despite promising progress within specific domains,  
33    they heavily rely on specialized data, techniques, and model architectures, suffering from additional  
34    resource and model design costs. More importantly, their focus on such specific aspects of visual  
35    coherence significantly limits their generalizability to a wider variety of applications. To this end, we  
36    propose a more challenging problem, ***Text-to-ImageSet (T2IS)*** generation, which aims to receive  
37    diverse user instructions and generate image sets that satisfy multifaceted consistency requirements.  
38    We address this challenge through two key contributions:

40    **(1) T2IS Benchmark:** We first present **T2IS-Bench**, the pioneering benchmark crafted specifically  
41    for T2IS generation. Derived from a comprehensive collection of real-world requirements, T2IS-  
42    Bench encompasses 596 representative tasks distributed across 26 detailed subcategories. Figure 1  
43    showcases illustrative examples of diverse tasks. Based on this, we introduce **T2IS-Eval**, an evaluation  
44    framework that automatically transforms user instructions into multifaceted assessment criteria across  
45    three aspects of consistency: ***identity***, ***style***, and ***logic***. Our framework then leverages the strong multi-  
46    image recognition capabilities of large-scale models [2] to serve as effective consistency evaluators  
47    and obtain logits of "Yes-or-No" answers for each criterion to assess the consistency fulfillment.  
48    More importantly, unlike previous methods that focus on a single aspect of consistency [14, 17],  
49    T2IS-Eval leverages transformed criteria to enable adaptive and interpretable consistency evaluation.

50    **(2) T2IS Generation:** We propose **AutoT2IS**, a training-free framework that maximally leverages the  
51    remarkable in-context generation capabilities of pretrained Diffusion Transformers (DiTs) [21, 24].  
52    Specifically, AutoT2IS first employs a structured recaptioning approach to systematically parse  
53    user instructions into informative prompts for both individual images and the global image set.  
54    Subsequently, we introduce a novel set-aware generation with the divide-and-conquer strategy: first  
55    binding individual images' prompts to their respective independent latents during early denoising  
56    stages to establish unique content characteristics, then integrating these latents through the multi-  
57    modal attention mechanism with a global consistency prompt. This enables each visual latent to  
58    simultaneously attend to its individual prompt, the global consistency prompt, and other images'  
59    visual latents. This way dynamically harmonizes visual elements to satisfy both image-level prompt  
60    alignment and set-level visual consistency, effectively bridging the gap between isolated image  
61    generation and set-level consistency without requiring specialized fine-tuning.

62    We conduct comprehensive experiments on the proposed T2IS-Bench to evaluate various approaches,  
63    including unified models (*e.g.*, Show-o [47], Janus-Pro [7]), compositional frameworks (*e.g.*, Gemini+Flux [24]), and agentic frameworks (*e.g.*, ISG-Agent [5], ChatDiT [22]). Our evaluation results  
64    indicate that while current methods have achieved excellent performance on prompt alignment,  
65    there remains a significant gap in achieving set-level consistency, highlighting the value of this  
66    underexplored research direction. Furthermore, both quantitative and qualitative results demonstrate  
67    that our proposed AutoT2IS substantially outperforms existing approaches, achieving significant  
68    improvements across three aspects of consistency. Notably, AutoT2IS also achieves competitive or  
69    even superior performance in specific domains compared to specialized methods.  
70

71    Interestingly, during the development of this work, advanced commercial image generation models  
72    like Gemini 2.0 [15] and GPT-4o [31] were released. This positions T2IS-Bench as the first platform  
73    to systematically evaluate and analyze their T2IS capabilities. Our assessments reveal that these  
74    models excel in consistency but often sacrifice image quality and prompt alignment. To this end, as  
75    shown in the bottom of Figure 1, AutoT2IS can be seamlessly integrated with models like GPT-4o

76 to further boost performance. By offering a transparent benchmark alongside a robust baseline, our  
77 work lays a strong foundation for the advancement of T2IS research.

## 78 2 Related Work

79 **Text-to-Image Generation.** With the advancement of diffusion models [18, 40], Text-to-Image  
80 (T2I) models [3, 4, 11, 12, 24, 31, 33, 35–37, 39] have demonstrated exceptional capabilities in  
81 generating high-quality images that accurately align with textual descriptions. Early approaches, such  
82 as DALL-E 2 [36] and Stable Diffusion [33], employed Latent Diffusion Models [37] (LDMs) with  
83 U-Net backbones [38] for efficient denoising in compressed latent space. Recent developments have  
84 advanced the field through the integration of transformer architectures into diffusion models, known  
85 as Diffusion Transformers (DiTs) [12, 24, 32]. These models leverage global attention mechanisms to  
86 capture complex dependencies in data, resulting in improved scalability and image fidelity. However,  
87 existing methods primarily focus on generating individual images, with limited exploration into the  
88 generation of coherent image sets. In this paper, we investigate and harness the in-context generation  
89 capabilities of DiTs to maximize their potential for addressing diverse T2IS tasks.

90 **Consistent Text-to-Image Generation.** Existing methods related to T2IS generation span various  
91 domains, primarily including consistent character generation [25, 29, 45] and storytelling [30, 51]. In  
92 character generation, training-based approaches such as PhotoMaker [25] aim to preserve identity  
93 via parameter-efficient fine-tuning, while training-free methods like One-Prompt-One-Story [29]  
94 exploit the contextual consistency of language models without additional training. In contrast to  
95 character generation, storytelling further demands style coherence across images. Make-a-Story [34]  
96 incorporates a visual memory module to capture contextual cues, and StoryDiffusion [51] introduces  
97 Consistent Self-Attention to enhance both identity and style consistency. However, when tackling  
98 new tasks like process generation [41] that require strong logical consistency for generating drawing  
99 processes, previous methods often fall short and need customized task-specific fine-tuning [41]. This  
100 suggests that consistencies in existing methods are generally domain-specific and lack generalizability.  
101 In this work, we propose T2IS generation to address a broad range of consistency requirements.

## 102 3 Methodology

103 Our goal is to develop T2IS methods that aim to generate coherent image sets from diverse user  
104 instructions. To achieve this objective, we first provide a comprehensive T2IS benchmark with an  
105 evaluation framework for diverse tasks in Sec. 3.1. In Sec. 3.2, we introduce our generation method,  
106 AutoT2IS, which consists of Structured Recaption and Set-Aware Generation to generate image sets  
107 that satisfy both image-level prompt alignment and set-level visual consistency.

### 108 3.1 T2IS Benchmark

109 **T2IS-Bench.** We first construct a comprehensive benchmark that reflects real-world demands on  
110 diverse consistency, as illustrated in Figure 2. The construction follows a multi-stage process:

111 *First*, we conduct an extensive analysis of real-world T2IS applications and categorize these into  
112 five major groups, covering diverse consistency requirements. For instance, Character Generation  
113 primarily focuses on identity preservation, while Process Generation emphasizes logical consistency.  
114 It is noteworthy that falling within the same groups does not mean the requirements of consistency  
115 are exactly the same. Rather, categories within each group vary in the extent to which they integrate  
116 and balance additional dimensions of consistency.

117 *Second*, we collect diverse data from two complementary sources. We incorporate established  
118 benchmarks including Consistency+ [29], ISG-Bench [5], IDEA-Bench [27] and others [21, 23,  
119 41]. Additionally, we curate examples from real-world platforms, such as e-commerce websites  
120 (*e.g.*, Amazon, Taobao), image generation communities (*e.g.*, Civitai), social media platforms (*e.g.*,  
121 RedNote), and other sources. From academic benchmarks, we select representative examples based  
122 on their diversity, while from platforms, we prioritize examples with user engagement and popularity.

123 *Third*, we standardize all these examples by employing MLLMs such as GPT-4o to transform these  
124 into conversational user instructions. We further augment the dataset through in-context learning  
125 techniques [10] to ensure balanced samples across all subcategories. All samples undergo rigorous

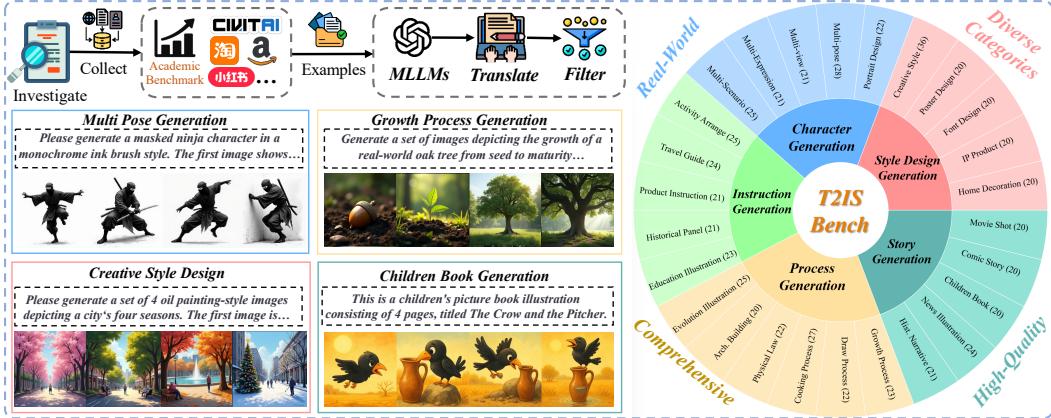


Figure 2: Illustration of the proposed T2IS-Bench. Upper Left: The collection process. Lower Left: Task examples displays. Right: Categories with their corresponding quantities.

126 verification by PhD-level domain experts, with duplicate or highly similar instances eliminated  
 127 through a combination of automated similarity detection using BertScore [50] and manual review.

128 Through the above process, we compile **T2IS-Bench**, containing 596 high-quality user instructions  
 129 distributed across 26 distinct subcategories. Our benchmark offers the most comprehensive coverage  
 130 to date, establishing a solid foundation for advancing T2IS. Detailed information regarding category  
 131 definitions, representative examples, and dataset distribution statistics is provided in Appendix A.

132 **T2IS-Eval.** We then propose **T2IS-Eval**, a comprehensive evaluation framework that assesses image  
 133 sets across three critical dimensions: *identity*, *style*, and *logical consistency*. For each dimension,  
 134 we employ LLMs to generate targeted evaluation questions as assessment criteria that reflect the  
 135 consistency requirements corresponding to the current user instructions. As illustrated in Figure 3,  
 136 the criteria across different dimensions include:

- 137 • **Identity** criteria measure the consistency of identity preservation across images. This includes  
 138 whether specific entities (*e.g.*, characters, objects) maintain consistent appearance features (*e.g.*,  
 139 size, color, shape), and whether their emotional expressions or key identifying elements persist.
- 140 • **Style** criteria assess the style uniform across images, including consistency in illustration style (*e.g.*,  
 141 watercolor, cartoon, 3D rendering), harmony in color palette usage, and others.
- 142 • **Logical** criteria evaluate whether the image set maintains reasonable causal relationships and  
 143 narrative coherence. This includes environmental consistency across scenes, proper depiction of  
 144 cause-and-effect relationships, and logical alignment between actions and their consequences.

145 Then, we systematically assess how well image sets satisfy these criteria. Inspired by VQAscore [28],  
 146 we discovered that large-scale MLLM [2] with strong multi-image recognition capabilities can  
 147 serve as excellent consistency judges. Based on this insight, we formulate our evaluation criteria as  
 148 questions that are sequentially applied to image pairs, obtaining logit scores for "Yes-or-No" answers  
 149 that quantify the degree of consistency fulfillment. In Appendix B, we demonstrate the effectiveness  
 150 of this approach and its advantages over direct scoring by MLLMs. Overall, unlike previous methods  
 151 that focus on a single consistency dimension [14], our approach leverages transformed criteria to  
 152 enables adaptive consistency evaluation with fine-grained measurements across multiple dimensions.

153 In addition to consistency evaluation, we also incorporate established methods to comprehensively  
 154 assess generation quality. Specifically, we employ MPS [49] to evaluate the aesthetic quality of gener-  
 155 ated images, and VQAscore [28] to measure the alignment between each image and its corresponding  
 156 instruction. Notably, all three evaluation components (aesthetics, prompt alignment, and visual  
 157 consistency) leverage vision-language models and employ a unified scoring mechanism based on  
 158 next-token prediction logits or text-image matching logits. Consequently, these quantitative metrics  
 159 are normalized to a consistent range [0, 1], providing better interpretability and holistic assessment.

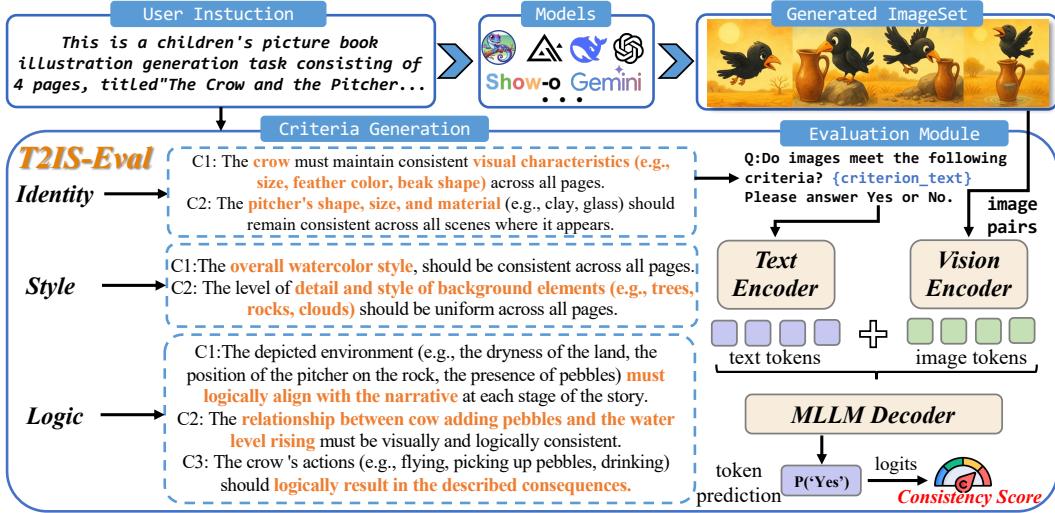


Figure 3: Overview of our T2IS-Eval. For each dimension, LLMs are prompted to generate 2-4 criteria; we only display a subset. Highlighted portions show the core elements for evaluation.

### 160 3.2 T2IS Generation

161 In this section, we introduce AutoT2IS, a training-free framework for T2IS generation. We aim to  
 162 maximally leverage the remarkable in-context generation capabilities [21] of pretrained Diffusion  
 163 Transformers (DiTs [24, 32]) to handle diverse T2IS tasks. We delineate two key steps, as depicted in  
 164 Figure 4. Specifically, AutoT2IS first utilizes a ***structured reception*** to parse user instructions into  
 165 informative prompts for both individual images and the complete image set, thereby capturing richer  
 166 semantic information. Subsequently, we introduce ***set-aware generation*** that effectively applies the  
 167 reciprocated individual prompts and consistency requirements to the entire image set, transforming  
 168 semantic contents into visually consistent outputs.

169 **Structured Reception.** Let  $y$  denote the given user instruction, which contains requirements for  
 170 individual images and overall consistency. We first employ LLMs to identify key content for each  
 171 image and consistency requirements for the entire set. The recognized elements are represented in a  
 172 structured format  $S = \{E, C\}$ , where  $E = \{e_1, e_2, \dots, e_n\}$  represents the content for  $n$  individual  
 173 images, and  $C = \{c_1, c_2, \dots, c_m\}$  denotes the  $m$  consistency requirements extracted from  $y$ .

174 Based on the structured information  $S$  and original user instruction  $y$ , we generate enhanced prompts  
 175 for each image and global consistency requirements:

$$p_i = f_{\text{recap}}(e_i, y) \quad \forall i \in \{1, 2, \dots, n\}, \quad g = f_{\text{consist}}(C, y) \quad (1)$$

176 where  $p_i$  represents the detailed prompt for the  $i$ -th image with denser fine-grained details to improve  
 177 fidelity, and  $g$  represents the global consistency description that ensures the model adheres to  
 178 overarching principles across all images. Functions  $f_{\text{recap}}$  and  $f_{\text{consist}}$  are implemented using LLMs to  
 179 expand the structured information into comprehensive textual descriptions.

180 **Set-Aware Generation.** While DiT models exhibit powerful in-context capabilities, existing ap-  
 181 proaches [20, 21, 42] often require fine-tuning to activate these abilities, which compromises the  
 182 model’s original generalization performance. As shown in Figure 4, we introduce set-aware genera-  
 183 tion that enables consistency requirements to be expressed through textual prompts and multi-modal  
 184 attention in a training-free way.

185 We empirically find that these few binding steps are sufficient to establish the unique content of  
 186 each image, allowing subsequent steps to focus on resolving consistency issues across the image  
 187 set. Inspired by this, we design a *divide-and-conquer* strategy to separate DiT’s native denoising  
 188 capabilities into two distinct phases. In the *divide* step, which is applied during the early stages of the  
 189 denoising process with a total timestep  $t$ , we generate independent noise latents  $x_t^i$  for each image  
 190  $i \in \{1, 2, \dots, n\}$  and denoise them separately with their corresponding prompts  $p_i$ . This binding  
 191 process is executed only during the initial  $r$  steps of the denoising process.

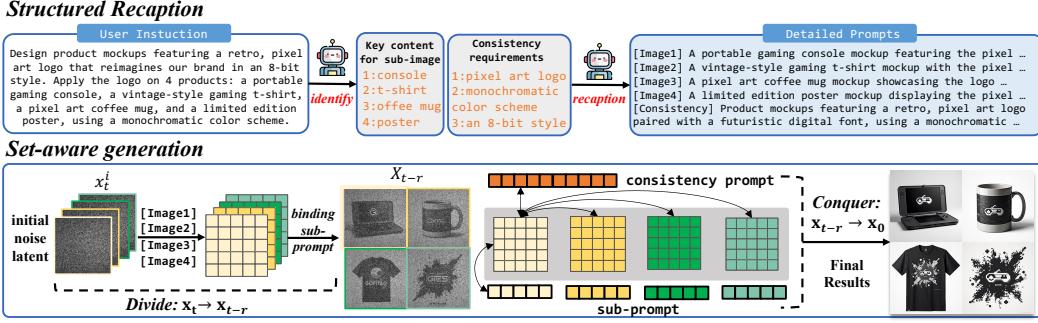


Figure 4: The pipeline of our AutoT2IS, consists of structured recaption and set-aware generation.

192 In the *conquer* phase, we concatenate all independently denoised latents  $x_{t-r}^i$  into a grid layout  
 193  $X_{t-r} = [x_{t-r}^1, x_{t-r}^2, \dots, x_{t-r}^n]$ . This concatenation enables previously isolated regions to collaborate  
 194 through DiT’s multi-modal attention mechanism [12, 32]. Formally, for each denoising step  $t' \in$   
 195  $\{t-r, t-r-1, \dots, 0\}$ , the attention computation can be expressed as:

$$\text{Attention}(Q, K, V, M) = \text{softmax} \left( \frac{QK^T \odot M}{\sqrt{d_k}} \right) V \quad (2)$$

196 where  $Q = W_Q X_{t'}$  represents queries derived from the concatenated latent  $X_{t'}$ ,  $K =$   
 197  $[W_K^{\text{text}} p; W_K^{\text{text}} g; W_K^{\text{image}} X_{t'}]$  and  $V = [W_V^{\text{text}} p; W_V^{\text{text}} g; W_V^{\text{image}} X_{t'}]$  index keys and values from both  
 198 the prompt embeddings  $P = [p_1, p_2, \dots, p_n, g]$  and the latent representation itself. The mask matrix  
 199  $M \in \mathbb{R}^{N \times (N_p + N_g + N)}$  is defined as:

$$M_{i,j} = \begin{cases} 1, & \text{if } i \in x_{t'}^k \text{ and } j \in \{p_k, g, X_{t'}\} \\ 0, & \text{otherwise} \end{cases}, \quad \text{where } k \in \{1, 2, \dots, n\} \quad (3)$$

200 This mask ensures that each visual latent from a specific region  $x_{t'}^k$  to simultaneously attend to three  
 201 critical elements, as illustrated in Figure 4:

- 202 • Its corresponding individual prompt  $p_i$  for precise semantic content preservation;  
 203 • The global consistency prompt  $g$  for harmonious coherence across the image set;  
 204 • Visual latents in other images for enhanced cross-image alignment and contextual integration.

205 Through this enhanced multi-modal attention mechanism, the model dynamically harmonizes visual  
 206 elements to satisfy both local alignment and global consistency requirements, effectively bridging the  
 207 gap between isolated image generation and set-level consistency without specialized fine-tuning.

208 **Extensions to image-conditioned and long-sequence generation.** Our method extends naturally  
 209 to image-conditioned and long-sequence T2IS generation. For image-conditioned tasks, we apply  
 210 inversion techniques [13, 46] to extract latent tokens from reference images, then implement token  
 211 replacement to preserve these conditioned regions while only denoising the remaining areas. This  
 212 maintains visual consistency with reference images. For long-sequence generation (more than 5  
 213 images), we employ a sliding window approach that generates new images based on previously  
 214 created windows, ensuring consistency across extended sequences while managing computational  
 215 constraints. Additional details are provided in Appendix C.

## 216 4 Experiments

217 **Implementation Details.** We utilize DeepSeek-R1 [9] as the foundational LLM for criteria generation  
 218 and structured recaptioning. Qwen2.5-VL-7B [2] serves as our consistency evaluator. For Set-Aware  
 219 generation, we employ FLUX.1-dev [24] as the base model, executing a total of 20 denoising steps,  
 220 i.e., 2 steps for the divide phase and the remaining for the conquer phase. For the concatenated latent,  
 221 we apply positional encoding at new positions to preserve the original multi-modal attention structure.

Table 1: Comparisons with various generalized generation models on T2IS-Bench. The upper shows open-source models, while the lower presents commercial models. The best and second-best results among the open-source models are highlighted in bold and underlined, respectively.

Model	Aesthetics	A Prompt Alignment			Visual Consistency			Avg.
		Entity	Attribute	Relation	Identity	Style	Logic	
Anole [8]	0.170	0.534	0.611	0.570	0.115	0.148	0.090	0.264
Show-o [47]	0.206	0.780	0.785	0.776	0.233	0.287	0.285	0.409
Janus-Pro [7]	0.333	0.787	0.785	0.781	0.224	0.272	0.300	0.435
Gemini & SD3 [12]	0.500	<b>0.796</b>	<b>0.794</b>	<b>0.789</b>	<u>0.287</u>	0.244	0.320	0.480
Gemini & Pixart [6]	0.447	0.743	0.765	0.747	0.206	0.279	0.268	0.440
Gemini & Hunyuan [26]	0.410	0.758	0.774	0.765	0.197	0.276	0.271	0.436
Gemini & Flux-1 [24]	<b>0.533</b>	<u>0.791</u>	<u>0.790</u>	<u>0.786</u>	0.249	0.302	<u>0.328</u>	<u>0.490</u>
ChatDit [22]	0.414	0.717	0.726	0.726	0.296	<u>0.326</u>	0.310	0.455
ISG-Agent [5]	0.256	0.667	0.703	0.682	0.146	0.178	0.163	0.338
<b>AutoT2IS (ours)</b>	<b>0.520</b>	0.729	0.756	0.743	<b>0.359</b>	<b>0.414</b>	<b>0.356</b>	<b>0.515</b>
Gemini 2.0 Flash [15]	0.430	0.738	0.747	0.743	0.428	0.383	0.392	0.509
GPT-4o [31]	0.445	0.663	0.683	0.693	0.400	0.463	0.383	0.501
<b>AutoT2IS + GPT-4o</b>	0.567	0.754	0.761	0.763	0.441	0.520	0.416	0.571

222 Due to space limitations, we present the most significant experimental results in the main text,  
223 including quantitative comparison of existing generalization and specialization methods with ours on  
224 T2IS-Bench (Sec. 4.1), qualitative visual comparisons to demonstrate the broad applications of our  
225 method (Sec. 4.2), and ablation studies on core components (Sec. 4.3). Please refer to the appendix  
226 for more comprehensive information. We provide detailed experimental settings in App.D, reliability  
227 analysis of T2IS-Eval in App.E, failure cases in App.G, limitations and future work in App.F.

## 228 4.1 Benchmarking T2IS: Quantitative Comparison of Generalized and Specialized Methods

229 **Experiment Setups.** We evaluate diverse frameworks for T2IS generation capabilities: (1) unified  
230 models like Show-o; (2) compositional models that use Gemini [43] to generate sub-captions with  
231 state-of-the-art T2I models; and (3) agentic models such as ChatDit that incorporate multi-step  
232 planning and generation processes. We also evaluate commercial models including Gemini 2.0 Flash  
233 and GPT-4o. Detailed setups for implementing these models are provided in Appendix H.

234 **Visual Consistency Challenges All Generalized Methods.** As illustrated in Table 1, most mod-  
235 els perform well on prompt alignment, with average scores all exceeding 0.5, demon-  
236 strating the significant progress T2I models have made in alignment capabilities. However, all models exhibit  
237 significant deficiencies in visual consistency. All open-source methods score below 0.35 across all  
238 consistency dimensions, highlighting the value of exploring this research direction. Comparatively,  
239 our method achieves the best consistency capabilities among open frameworks, showing notable  
240 improvements over previous methods and delivering the best overall performance.

241 **Commercial Models Still Need Improvement.** As shown in the bottom of Table 1, we evaluate  
242 state-of-the-art commercial models [15, 31]. These models, benefiting from large-scale pretraining,  
243 demonstrate superior instruction understanding and better consistency than open-source alternatives.  
244 However, our results reveal that they still struggle with T2IS generation tasks. During T2IS genera-  
245 tion, their aesthetic quality and prompt alignment significantly deteriorate, even performing worse  
246 than many open-source models. This shortcoming can be mitigated by incorporating Auto-T2IS’s  
247 structured recaption to enrich semantic information, which achieves significant improvements in  
248 quality, alignment, and consistency. These findings underscore the necessity for improved T2IS task  
249 instruction comprehension in commercial models.

250 **AutoT2IS Matches/Exceeds Specialized Methods.** Table 2 provides comparisons between our  
251 unified generation approach and specialized methods. These specialized methods are typically  
252 tailored and optimized for specific domains, allowing them to perform exceptionally well in those  
253 areas but limiting their adaptability across a wide range of T2IS tasks. The results in their respective  
254 domains on T2IS show that our method achieves competitive or even superior performance in areas  
255 where specialized methods have been pre-trained. This suggests that AutoT2IS effectively explores

Table 2: Comparisons of specialized methods and our unified workflow across different T2IS domains.

Model	Aesthetics	A Prompt Alignment			Visual Consistency			Avg.
		Entity	Attribute	Relation	Identity	Style	Logic	
<b>Character Generation: Multi-view, Multi-scenario, Portrait Design</b>								
PhotoMaker [25]	0.413	0.793	0.799	0.762	0.482	0.458	0.450	0.552
OnePrompt [29]	0.600	0.805	0.807	0.766	0.520	0.512	0.470	0.590
X-Flux [1]	0.428	0.834	0.804	0.780	0.531	0.529	0.521	0.553
<b>AutoT2IS (ours)</b>	<b>0.557</b>	<b>0.798</b>	<b>0.788</b>	<b>0.764</b>	<b>0.609</b>	<b>0.520</b>	<b>0.509</b>	<b>0.619</b>
<b>Style Design Generation: Creative Style, Font, IP product design</b>								
IPAdapter-Flux [44]	0.514	0.856	0.833	0.846	0.386	0.455	0.517	<b>0.582</b>
<b>AutoT2IS (ours)</b>	0.475	0.789	0.791	0.792	0.339	0.391	0.478	0.533
<b>Story Generation: Movie Shot, Comic Story, Children Book</b>								
ICLora-Story [21]	0.408	0.711	0.739	0.702	0.078	0.190	0.119	0.361
StoryDiffusion [51]	0.581	0.684	0.728	0.670	0.182	0.207	0.141	0.413
Story-Adapter [30]	0.578	0.749	0.766	0.734	0.181	0.328	0.226	<b>0.463</b>
<b>AutoT2IS (ours)</b>	0.534	0.629	0.698	0.660	0.262	0.405	0.237	0.456
<b>Process Generation: Growth, Draw, Building</b>								
MakeAnything [41]	0.358	0.636	0.665	0.652	0.329	0.340	0.222	0.408
<b>AutoT2IS (ours)</b>	<b>0.578</b>	<b>0.727</b>	<b>0.754</b>	<b>0.739</b>	<b>0.314</b>	<b>0.346</b>	<b>0.274</b>	<b>0.493</b>

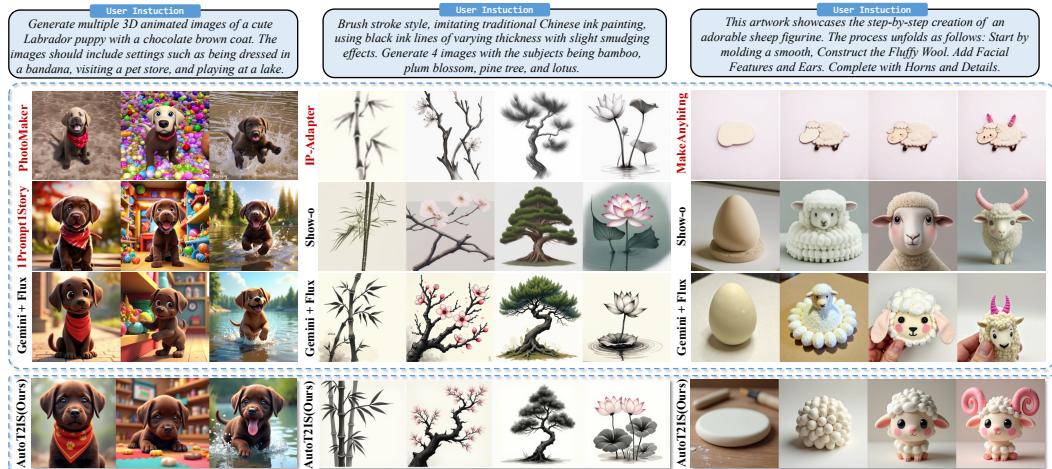


Figure 5: Qualitative comparison between our AutoT2IS and existing methods on common T2IS tasks. Black labels indicate generalized methods, while red labels represent specialized methods.

256 broader consistency principles while avoiding overfitting, thereby maintaining the model’s inherent  
257 generative capabilities. These findings highlight the value of unified methods for diverse T2IS tasks.

## 258 4.2 Visualizing T2IS: Qualitative Comparison and Diverse Application Showcase

259 **AutoT2IS Excels at Common T2IS Tasks.** Figure 5 presents a qualitative comparison between our  
260 AutoT2IS and existing methods across several common T2IS tasks. As demonstrated, generalized  
261 methods like Show-o and Flux exhibit significant limitations in identity preservation, style uniform,  
262 and logical coherence. In contrast, AutoT2IS produces more visually consistent image sets. When  
263 evaluated against specialized methods in multi-scenario character generation, our approach main-  
264 tains exceptional identity consistency while ensuring precise prompt alignment, achieving results  
265 comparable to the SOTA method 1Prompt1Story [29]. Notably, when it comes to challenging logical  
266 consistency generation tasks, existing methods like MakeAnything [41], despite being specifically  
267 trained for these scenarios, often compromise their fundamental generation capabilities, resulting  
268 in diminished image quality. Conversely, our training-free method preserves the model’s original  
269 generation capabilities and simultaneously maximizing its inherent consistency potential.



Figure 6: AutoT2IS enables novel T2IS applications with significant real-world value.

270 **AutoT2IS Enables More Real-world Applications.** Beyond common T2IS tasks, our AutoT2IS en-  
 271 ables to address previously underexplored applications with significant real-world value, as illustrated  
 272 in Figure 6. These applications have been largely overlooked due to the lack of a unified approach to  
 273 varying degrees of different consistency. For instance, IP Product Design (Task 2) requires coherence  
 274 between IP concept identity and product style. Similarly, Growth Process Generation and Physical  
 275 Law Illustration (Tasks 4 and 6) demand comprehensive consistency across identity, style, and logical  
 276 dimensions. The diversity of these tasks demonstrates that focusing on a single aspect of consistency  
 277 severely limits the range of potential applications, highlighting the value of unified methods for T2IS.

### 278 4.3 Ablation Studies

279 To better understand the contribution of each stage in AutoT2IS, we conduct ablation studies. Table 3  
 280 presents the quantitative results across different evaluation dimensions.

281 **Importance of Structured Recaption (SR).** The *structured recaption* plays a critical role  
 282 by reformulating user inputs into more struc-  
 283 tured and detailed descriptions that better guide  
 284 the image generation process. When we re-  
 285 place the *structured recaption* with a baseline  
 286 approach that simply uses an LLM to gen-  
 287 erate sub-captions, we observe a significant drop  
 288 across all evaluation metrics. These results high-  
 289 light the importance of semantic completeness and relational coherence for effective T2IS generation.

290 **Importance of Set-Aware Generation (SG).** The set-aware generation is responsible for creating a  
 291 consistent set across multiple images, while ensuring the quality and alignment of each sub-image. We  
 292 replace this stage with a baseline approach that directly concatenates prompts to guide the model in  
 293 generating multi-grid images [21, 22]. As shown in results, despite achieving reasonable consistency,  
 294 the quality of individual images and text-image alignment significantly decreased. This demonstrates  
 295 the importance of the divide-and-conquer approach in Set-Aware Generation, where the model first  
 296 understands the content of individual image before rendering them with visual consistency.

297 We include additional ablation with visualizations in Appendix I for more intuitive understanding.

## 299 5 Conclusion

300 In this paper, we advance the field of Text-to-ImageSet (T2IS) by introducing innovative frameworks  
 301 for both T2IS generation and evaluation. We present T2IS-Bench, a comprehensive benchmark  
 302 featuring 596 representative tasks across 26 detailed subcategories. Alongside this, we introduce  
 303 T2IS-Eval, an evaluation framework that enables diverse consistency assessment. Furthermore,  
 304 we introduce AutoT2IS, a training-free framework that fully exploits the exceptional in-context  
 305 generation capabilities of DiT models, effectively harmonizing visual elements to meet both image-  
 306 level prompt alignment and set-level visual consistency. By providing a transparent benchmark and a  
 307 robust baseline, our work establishes a solid foundation for the future of T2IS research.

308 **References**

- 309 [1] XLabs AI. x-flux: Flux model fine-tuning scripts. <https://github.com/XLabs-AI/x-flux>,  
310 2024. Apache License 2.0.
- 311 [2] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang,  
312 Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang  
313 Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen  
314 Cheng, Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. Qwen2.5-v1 technical report.  
315 *arXiv preprint arXiv:2502.13923*, 2025.
- 316 [3] Jason Baldridge, Kevin Shih, Yael Li, Zongyi Wang, Shweta Prabhudesai, Sashi Cheluvaraju,  
317 Xin Li, Lucy Chai, Miao Ding, Bryan Catanzaro, et al. Imagen 2: Tuning text-to-image diffusion  
318 models for photorealism and generalization. *arXiv preprint arXiv:2401.18680*, 2024.
- 319 [4] James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jianfeng Wang, Linjie Li, Long Ouyang,  
320 Juntang Zhuang, Joyce Lee, Yufei Guo, et al. Improving image generation with better captions.  
321 *Computer Science*. <https://cdn.openai.com/papers/dall-e-3.pdf>, 2(3):8, 2023.
- 322 [5] Dongping Chen, Ruoxi Chen, Shu Pu, Zhaoyi Liu, Yanru Wu, Caixi Chen, Benlin Liu, Yue  
323 Huang, Yao Wan, Pan Zhou, and Ranjay Krishna. Interleaved scene graph for interleaved  
324 text-and-image generation assessment. *arXiv preprint arXiv:2411.17188*, 2024.
- 325 [6] Junsong Chen, Jincheng Yu, Chongjian Ge, Lewei Yao, Enze Xie, Yue Wu, Zhongdao Wang,  
326 James Kwok, Ping Luo, Huchuan Lu, and Zhenguo Li. Pixart- $\alpha$ : Fast training of diffusion  
327 transformer for photorealistic text-to-image synthesis, 2023.
- 328 [7] Xiaokang Chen, Zhiyu Wu, Xingchao Liu, Zizheng Pan, Wen Liu, Zhenda Xie, Xingkai Yu,  
329 and Chong Ruan. Janus-pro: Unified multimodal understanding and generation with data and  
330 model scaling. *arXiv preprint arXiv:2501.17811*, 2025.
- 331 [8] Ethan Chern, Jiadi Su, Yan Ma, and Pengfei Liu. Anole: An open, autoregressive, native large  
332 multimodal models for interleaved image-text generation. *arXiv preprint arXiv:2407.06135*,  
333 2024.
- 334 [9] DeepSeek-AI. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement  
335 learning, 2025. URL <https://arxiv.org/abs/2501.12948>.
- 336 [10] Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan Ma, Rui Li, Heming Xia, Jingjing  
337 Xu, Zhiyong Wu, Tianyu Liu, et al. A survey on in-context learning. *arXiv preprint  
338 arXiv:2301.00234*, 2022.
- 339 [11] Patrick Esser, Robin Rombach, and Björn Ommer. Taming transformers for high-resolution  
340 image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern  
341 Recognition*, pages 12873–12883, 2021.
- 342 [12] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini,  
343 Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transform-  
344 ers for high-resolution image synthesis. In *Forty-first International Conference on Machine  
345 Learning*, 2024.
- 346 [13] Haoran Feng, Zehuan Huang, Lin Li, Hairong Lv, and Lu Sheng. Personalize anything for free  
347 with diffusion transformer. *arXiv preprint arXiv:2503.12590*, 2025.
- 348 [14] Stephanie Fu, Netanel Tamir, Shobhit Sundaram, Lucy Chai, Richard Zhang, Tali Dekel, and  
349 Phillip Isola. Dreamsim:learning new dimensions of human visual similarity using synthetic  
350 data. In *Advances in Neural Information Processing Systems*, volume 36, pages 50742–50768,  
351 2023.
- 352 [15] Google Cloud. Generate images with gemini | generative ai on vertex ai, 2025. URL  
353 [https://cloud.google.com/vertex-ai/generative-ai/docs/multimodal/  
354 image-generation](https://cloud.google.com/vertex-ai/generative-ai/docs/multimodal/image-generation). Accessed: 2025-05-09.

- 355 [16] Amir Hertz, Andrey Voynov, Shlomi Fruchter, and Daniel Cohen-Or. Style aligned image  
356 generation via shared attention. In *Proceedings of the IEEE/CVF Conference on Computer*  
357 *Vision and Pattern Recognition*, pages 4775–4785, 2024.
- 358 [17] Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore:  
359 A reference-free evaluation metric for image captioning. In *Empirical Methods in Natural*  
360 *Language Processing (EMNLP)*, 2021.
- 361 [18] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances*  
362 *in neural information processing systems*, 33:6840–6851, 2020.
- 363 [19] Jiehui Huang, Xiao Dong, Wenhui Song, Hanhui Li, Jun Zhou, Yuhao Cheng, Shutao Liao, Long  
364 Chen, Yiqiang Yan, Shengcai Liao, et al. Consistentid: Portrait generation with multimodal  
365 fine-grained identity preserving. *arXiv preprint arXiv:2404.16771*, 2024.
- 366 [20] Lianghua Huang, Wei Wang, Zhi-Fan Wu, Huanzhang Dou, Yupeng Shi, Yutong Feng, Chen  
367 Liang, Yu Liu, and Jingren Zhou. Group diffusion transformers are unsupervised multitask  
368 learners. 2024.
- 369 [21] Lianghua Huang, Wei Wang, Zhi-Fan Wu, Yupeng Shi, Huanzhang Dou, Chen Liang, Yutong  
370 Feng, Yu Liu, and Jingren Zhou. In-context lora for diffusion transformers. *arXiv preprint*  
371 *arXiv:2410.23775*, 2024.
- 372 [22] Lianghua Huang, Wei Wang, Zhi-Fan Wu, Yupeng Shi, Chen Liang, Tong Shen, Han Zhang,  
373 Huanzhang Dou, Yu Liu, and Jingren Zhou. Chatdit: A training-free baseline for task-agnostic  
374 free-form chatting with diffusion transformers. 2024.
- 375 [23] Zehuan Huang, Yuan-Chen Guo, Haoran Wang, Ran Yi, Lizhuang Ma, Yan-Pei Cao, and  
376 Lu Sheng. Mv-adapter: Multi-view consistent image generation made easy. *arXiv preprint*  
377 *arXiv:2412.03632*, 2024.
- 378 [24] Black Forest Labs. Flux: Inference repository. <https://github.com/black-forest-labs/flux>, 2024. Accessed: 2024-10-25.
- 380 [25] Zhen Li, Mingdeng Cao, Xintao Wang, Zhongang Qi, Ming-Ming Cheng, and Ying Shan.  
381 Photomaker: Customizing realistic human photos via stacked id embedding. In *IEEE Conference*  
382 *on Computer Vision and Pattern Recognition (CVPR)*, 2024.
- 383 [26] Zhimin Li, Jianwei Zhang, Qin Lin, Jiangfeng Xiong, Yanxin Long, Xinchi Deng, Yingfang  
384 Zhang, Xingchao Liu, Minbin Huang, Zedong Xiao, Dayou Chen, Jiajun He, Jiahao Li, Wenyue  
385 Li, Chen Zhang, Rongwei Quan, Jianxiang Lu, Jiabin Huang, Xiaoyan Yuan, Xiaoxiao Zheng,  
386 Yixuan Li, Jihong Zhang, Chao Zhang, Meng Chen, Jie Liu, Zheng Fang, Weiyan Wang, Jinbao  
387 Xue, Yangyu Tao, Jianchen Zhu, Kai Liu, Sihuan Lin, Yifu Sun, Yun Li, Dongdong Wang,  
388 Mingtao Chen, Zhichao Hu, Xiao Xiao, Yan Chen, Yuhong Liu, Wei Liu, Di Wang, Yong Yang,  
389 Jie Jiang, and Qinglin Lu. Hunyuan-dit: A powerful multi-resolution diffusion transformer with  
390 fine-grained chinese understanding, 2024.
- 391 [27] Chen Liang, Lianghua Huang, Jingwu Fang, Huanzhang Dou, Wei Wang, Zhi-Fan Wu, Yupeng  
392 Shi, Junge Zhang, Xin Zhao, and Yu Liu. Idea-bench: How far are generative models from  
393 professional designing? *arXiv preprint arXiv:2412.11767*, 2024.
- 394 [28] Zhiqiu Lin, Deepak Pathak, Baiqi Li, Jiayao Li, Xide Xia, Graham Neubig, Pengchuan Zhang,  
395 and Deva Ramanan. Evaluating text-to-visual generation with image-to-text generation. In  
396 *European Conference on Computer Vision*, pages 366–384. Springer, 2024.
- 397 [29] Tao Liu, Kai Wang, Senmao Li, Joost van de Weijer, Fahad Shahbaz Khan, Shiqi Yang, Yaxing  
398 Wang, Jian Yang, and Ming-Ming Cheng. One-prompt-one-story: Free-lunch consistent text-  
399 to-image generation using a single prompt. In *The Thirteenth International Conference on*  
400 *Learning Representations*, 2025. URL <https://openreview.net/forum?id=cD1k12QKv1>.
- 401 [30] Jiawei Mao, Xiaoke Huang, Yunfei Xie, Yuanqi Chang, Mude Hui, Bingjie Xu, and Yuyin Zhou.  
402 Story-Adapter: A Training-free Iterative Framework for Long Story Visualization, 2024.

- 403 [31] OpenAI. Introducing 4o image generation. <https://openai.com/index/introducing-4o-image-generation/>, 2025. Accessed: 2025-03-25.
- 404
- 405 [32] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings*  
406 *of the IEEE/CVF International Conference on Computer Vision*, pages 4195–4205, 2023.
- 407 [33] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe  
408 Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image  
409 synthesis. *arXiv preprint arXiv:2307.01952*, 2023.
- 410 [34] Tanzila Rahman, Hsin-Ying Lee, Jian Ren, Sergey Tulyakov, Shweta Mahajan, and Leonid  
411 Sigal. Make-a-story: Visual memory conditioned consistent story generation. In *Proceedings of*  
412 *the IEEE/CVF conference on computer vision and pattern recognition*, pages 2493–2502, 2023.
- 413 [35] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark  
414 Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on*  
415 *Machine Learning*, pages 8821–8831. PMLR, 2021.
- 416 [36] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical  
417 text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.
- 418 [37] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
419 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF*  
420 *conference on computer vision and pattern recognition*, pages 10684–10695, 2022.
- 421 [38] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks  
422 for biomedical image segmentation. In *Medical image computing and computer-assisted*  
423 *intervention–MICCAI 2015: 18th international conference, Munich, Germany, October 5–9,*  
424 *2015, proceedings, part III 18*, pages 234–241. Springer, 2015.
- 425 [39] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Tim  
426 Salimans, Jonathan Ho, David J Fleet, Phillip Isola, et al. Photorealistic text-to-image diffusion  
427 models with deep language understanding. *arXiv preprint arXiv:2205.11487*, 2022.
- 428 [40] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv*  
429 *preprint arXiv:2010.02502*, 2020.
- 430 [41] Yiren Song, Cheng Liu, and Mike Zheng Shou. Makeanything: Harnessing diffusion trans-  
431 formers for multi-domain procedural sequence generation. *arXiv preprint arXiv:2502.01572*,  
432 2025.
- 433 [42] Zhenxiong Tan, Songhua Liu, Xingyi Yang, Qiaochu Xue, and Xinchao Wang. Ominicontrol:  
434 Minimal and universal control for diffusion transformer. *arXiv preprint arXiv:2411.15098*,  
435 2024.
- 436 [43] Gemini Team, Rohan Anil, Sébastien Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut,  
437 Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: a family of highly  
438 capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- 439 [44] InstantX Team. Instantx flux.1-dev ip-adapter page, 2024.
- 440 [45] Yoad Tewel, Omri Kaduri, Rinon Gal, Yoni Kasten, Lior Wolf, Gal Chechik, and Yuval Atzmon.  
441 Training-free consistent text-to-image generation. *ACM Transactions on Graphics (TOG)*, 43  
442 (4):52, 2024. doi: 10.1145/3658157. URL <https://doi.org/10.1145/3658157>.
- 443 [46] Jiangshan Wang, Junfu Pu, Zhongang Qi, Jiayi Guo, Yue Ma, Nisha Huang, Yuxin Chen,  
444 Xiu Li, and Ying Shan. Taming rectified flow for inversion and editing. *arXiv preprint*  
445 *arXiv:2411.04746*, 2024.
- 446 [47] Jinheng Xie, Weijia Mao, Zechen Bai, David Junhao Zhang, Weihao Wang, Kevin Qinghong  
447 Lin, Yuchao Gu, Zhijie Chen, Zhenheng Yang, and Mike Zheng Shou. Show-o: One single trans-  
448 former to unify multimodal understanding and generation. *arXiv preprint arXiv:2408.12528*,  
449 2024.

- 450 [48] Hu Ye, Jun Zhang, Sibo Liu, Xiao Han, and Wei Yang. Ip-adapter: Text compatible image  
451 prompt adapter for text-to-image diffusion models. 2023.
- 452 [49] Sixian Zhang, Bohan Wang, Junqiang Wu, Yan Li, Tingting Gao, Di Zhang, and Zhongyuan  
453 Wang. Learning multi-dimensional human preference for text-to-image generation. In *Pro-*  
454 *ceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages  
455 8018–8027, 2024.
- 456 [50] Tianyi Zhang\*, Varsha Kishore\*, Felix Wu\*, Kilian Q. Weinberger, and Yoav Artzi. Bertscore:  
457 Evaluating text generation with bert. In *International Conference on Learning Representations*,  
458 2020. URL <https://openreview.net/forum?id=SkeHuCVFDr>.
- 459 [51] Yupeng Zhou, Daquan Zhou, Ming-Ming Cheng, Jiashi Feng, and Qibin Hou. Storydiffusion:  
460 Consistent self-attention for long-range image and video generation. *Advances in Neural*  
461 *Information Processing Systems*, 37:110315–110340, 2024.

462 **NeurIPS Paper Checklist**

463 **1. Claims**

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

466 Answer: [Yes]

467 Justification: The abstract and introduction clearly and accurately reflect our core contribu-  
468 tions and scope. We introduce a new problem setting (Text-to-ImageSet generation),  
469 constructs a benchmark dataset (T2IS-Bench), proposes a corresponding evaluation frame-  
470 work (T2IS-Eval), and presents a training-free solution (AutoT2IS). The claims are supported  
471 by experimental results demonstrating the superiority of AutoT2IS over existing approaches  
472 and its generalization across a broad set of real-world applications.

473 Guidelines:

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

483 **2. Limitations**

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

485 Answer: [Yes]

486 Justification: The detailed limitation analyses are available in the Appendix F.

487 Guidelines:

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

514 **3. Theory assumptions and proofs**

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

517 Answer: [NA]

518 Justification: The paper does not contain any theoretical results, assumptions, or formal  
519 proofs. Its contributions are empirical and methodological in nature, focusing on benchmark  
520 construction, evaluation framework design, and a training-free generative method.

521 Guidelines:

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

#### 532 4. Experimental result reproducibility

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

536 Answer: [Yes]

537 Justification: The detailed experimental settings are available in the Experiments section 4  
538 of the main paper and the Appendix D. All data and code will be made public.

539 Guidelines:

- 540 • The answer NA means that the paper does not include experiments.
- 541 • If the paper includes experiments, a No answer to this question will not be perceived  
542 well by the reviewers: Making the paper reproducible is important, regardless of  
543 whether the code and data are provided or not.
- 544 • If the contribution is a dataset and/or model, the authors should describe the steps taken  
545 to make their results reproducible or verifiable.
- 546 • Depending on the contribution, reproducibility can be accomplished in various ways.  
547 For example, if the contribution is a novel architecture, describing the architecture fully  
548 might suffice, or if the contribution is a specific model and empirical evaluation, it may  
549 be necessary to either make it possible for others to replicate the model with the same  
550 dataset, or provide access to the model. In general, releasing code and data is often  
551 one good way to accomplish this, but reproducibility can also be provided via detailed  
552 instructions for how to replicate the results, access to a hosted model (e.g., in the case  
553 of a large language model), releasing of a model checkpoint, or other means that are  
554 appropriate to the research performed.
- 555 • While NeurIPS does not require releasing code, the conference does require all submis-  
556 sions to provide some reasonable avenue for reproducibility, which may depend on the  
557 nature of the contribution. For example
  - 558 (a) If the contribution is primarily a new algorithm, the paper should make it clear how  
559 to reproduce that algorithm.
  - 560 (b) If the contribution is primarily a new model architecture, the paper should describe  
561 the architecture clearly and fully.
  - 562 (c) If the contribution is a new model (e.g., a large language model), then there should  
563 either be a way to access this model for reproducing the results or a way to reproduce  
564 the model (e.g., with an open-source dataset or instructions for how to construct  
565 the dataset).
  - 566 (d) We recognize that reproducibility may be tricky in some cases, in which case  
567 authors are welcome to describe the particular way they provide for reproducibility.  
568 In the case of closed-source models, it may be that access to the model is limited in

569 some way (e.g., to registered users), but it should be possible for other researchers  
570 to have some path to reproducing or verifying the results.

571 **5. Open access to data and code**

572 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
573 tions to faithfully reproduce the main experimental results, as described in supplemental  
574 material?

575 Answer: [Yes]

576 Justification: We provide the details of benchmark data in Appendix A, and detailed experi-  
577 mental settings in Appendix D. All data and code will be made public upon acceptance.

578 Guidelines:

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

598 **6. Experimental setting/details**

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

602 Answer: [Yes]

603 Justification: The detailed experimental settings are available in the Experiments section 4  
604 of the main paper and the Appendix D.

605 Guidelines:

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

611 **7. Experiment statistical significance**

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

614 Answer: [Yes]

615 Justification: We conduct multiple runs to ensure the reliability of the evaluation.

616 Guidelines:

- 617 • The answer NA means that the paper does not include experiments.
- 618 • The authors should answer "Yes" if the results are accompanied by error bars, confi-  
619 dence intervals, or statistical significance tests, at least for the experiments that support  
620 the main claims of the paper.

- 621           • The factors of variability that the error bars are capturing should be clearly stated (for  
622           example, train/test split, initialization, random drawing of some parameter, or overall  
623           run with given experimental conditions).  
624           • The method for calculating the error bars should be explained (closed form formula,  
625           call to a library function, bootstrap, etc.)  
626           • The assumptions made should be given (e.g., Normally distributed errors).  
627           • It should be clear whether the error bar is the standard deviation or the standard error  
628           of the mean.  
629           • It is OK to report 1-sigma error bars, but one should state it. The authors should  
630           preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis  
631           of Normality of errors is not verified.  
632           • For asymmetric distributions, the authors should be careful not to show in tables or  
633           figures symmetric error bars that would yield results that are out of range (e.g. negative  
634           error rates).  
635           • If error bars are reported in tables or plots, The authors should explain in the text how  
636           they were calculated and reference the corresponding figures or tables in the text.

637           **8. Experiments compute resources**

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

641           Answer: [Yes]

642           Justification: All experiments are performed using NVIDIA L40s GPUs. The detailed  
643           experimental settings are available in the Experiments section 4 of the main paper and the  
644           Appendix D.

645           Guidelines:

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

654           **9. Code of ethics**

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

657           Answer: [Yes]

658           Justification: The research conforms to the NeurIPS Code of Ethics. All datasets used in this  
659           study are publicly available and widely used in the community. No personally identifiable  
660           information or sensitive content is involved.

661           Guidelines:

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

667           **10. Broader impacts**

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

670           Answer: [Yes]

671 Justification: The paper discusses both the positive and potential negative societal impacts.  
672 On the positive side, the proposed T2IS task and framework could significantly improve  
673 the controllability and practical utility of text-to-image generation models, enabling more  
674 coherent and useful image generation for applications such as education, creative design,  
675 e-commerce, and assistive tools. However, the work also raises potential concerns: gener-  
676 ating consistent image sets based on user instructions could be misused to mass-produce  
677 coordinated synthetic images for misinformation, propaganda, or impersonation purposes.  
678 To mitigate these risks, we suggest responsible release of model checkpoints and emphasize  
679 the importance of watermarking, usage monitoring, and ethical deployment practices.

680 Guidelines:

- 681 • The answer NA means that there is no societal impact of the work performed.
- 682 • If the authors answer NA or No, they should explain why their work has no societal  
683 impact or why the paper does not address societal impact.
- 684 • Examples of negative societal impacts include potential malicious or unintended uses  
685 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations  
686 (e.g., deployment of technologies that could make decisions that unfairly impact specific  
687 groups), privacy considerations, and security considerations.
- 688 • The conference expects that many papers will be foundational research and not tied  
689 to particular applications, let alone deployments. However, if there is a direct path to  
690 any negative applications, the authors should point it out. For example, it is legitimate  
691 to point out that an improvement in the quality of generative models could be used to  
692 generate deepfakes for disinformation. On the other hand, it is not needed to point out  
693 that a generic algorithm for optimizing neural networks could enable people to train  
694 models that generate Deepfakes faster.
- 695 • The authors should consider possible harms that could arise when the technology is  
696 being used as intended and functioning correctly, harms that could arise when the  
697 technology is being used as intended but gives incorrect results, and harms following  
698 from (intentional or unintentional) misuse of the technology.
- 699 • If there are negative societal impacts, the authors could also discuss possible mitigation  
700 strategies (e.g., gated release of models, providing defenses in addition to attacks,  
701 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from  
702 feedback over time, improving the efficiency and accessibility of ML).

## 703 11. Safeguards

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

707 Answer: [Yes]

708 Justification: The paper acknowledges the potential risks of misuse associated with releasing  
709 image generation frameworks and datasets. To mitigate these concerns, the released T2IS-  
710 Bench dataset has undergone a manual filtering process to ensure that it does not contain  
711 unsafe, offensive, or privacy-violating content. For the AutoT2IS model code, we provide  
712 detailed usage documentation and include explicit terms of use to discourage malicious  
713 applications such as mass-produced disinformation. In addition, we do not release fine-  
714 tuned weights or specialized prompt templates that could be easily misused without further  
715 safeguards. We encourage responsible research use and welcome community feedback to  
716 further improve safety practices.

717 Guidelines:

- 718 • The answer NA means that the paper poses no such risks.
- 719 • Released models that have a high risk for misuse or dual-use should be released with  
720 necessary safeguards to allow for controlled use of the model, for example by requiring  
721 that users adhere to usage guidelines or restrictions to access the model or implementing  
722 safety filters.
- 723 • Datasets that have been scraped from the Internet could pose safety risks. The authors  
724 should describe how they avoided releasing unsafe images.

- 725           • We recognize that providing effective safeguards is challenging, and many papers do  
726           not require this, but we encourage authors to take this into account and make a best  
727           faith effort.

728           **12. Licenses for existing assets**

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

732           Answer: [Yes]

733           Justification: The paper makes use of publicly available pretrained Diffusion Transformers  
734           and other generative models, such as Stable Diffusion, which are properly cited in the paper.  
735           We explicitly mention the versions used and adhere to their respective licenses—for instance,  
736           Stable Diffusion under the CreativeML Open RAIL-M license. For any datasets used for  
737           evaluation or benchmarking, we ensure their original creators are credited and that their  
738           terms of use are respected. All reused assets are only employed within the allowed scope of  
739           their licenses.

740           Guidelines:

- 741           • The answer NA means that the paper does not use existing assets.
- 742           • The authors should cite the original paper that produced the code package or dataset.
- 743           • The authors should state which version of the asset is used and, if possible, include a  
744            URL.
- 745           • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 746           • For scraped data from a particular source (e.g., website), the copyright and terms of  
747            service of that source should be provided.
- 748           • If assets are released, the license, copyright information, and terms of use in the  
749            package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets)  
750            has curated licenses for some datasets. Their licensing guide can help determine the  
751            license of a dataset.
- 752           • For existing datasets that are re-packaged, both the original license and the license of  
753            the derived asset (if it has changed) should be provided.
- 754           • If this information is not available online, the authors are encouraged to reach out to  
755            the asset's creators.

756           **13. New assets**

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

759           Answer: [Yes]

760           Justification: The paper introduces several new assets, including the T2IS-Bench dataset, the  
761           T2IS-Eval evaluation toolkit, and the AutoT2IS framework. All assets are well documented  
762           and will be released with accompanying instructions on usage, data format, evaluation  
763           protocols, and limitations. The dataset construction process is described in detail in the  
764           paper, and any content included in T2IS-Bench was either synthetically generated or curated  
765           from publicly available and appropriately licensed sources, avoiding privacy violations or  
766           copyrighted material. At submission time, an anonymized version of the assets is provided  
767           to ensure compliance with double-blind review requirements. Upon publication, all assets  
768           will be released under a permissive license with appropriate citation and usage guidelines.

769           Guidelines:

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

778     **14. Crowdsourcing and research with human subjects**

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

782     Answer: [NA]

783     Justification: This work does not involve crowdsourcing or research with human subjects.

784     Guidelines:

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

793     **15. Institutional review board (IRB) approvals or equivalent for research with human  
794         subjects**

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

799     Answer: [NA]

800     Justification: This work does not involve human subjects or crowdsourcing experiments and  
801     therefore does not require IRB or equivalent ethical approval.

802     Guidelines:

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

813     **16. Declaration of LLM usage**

814     Question: Does the paper describe the usage of LLMs if it is an important, original, or  
815     non-standard component of the core methods in this research? Note that if the LLM is used  
816     only for writing, editing, or formatting purposes and does not impact the core methodology,  
817     scientific rigorousness, or originality of the research, declaration is not required.

818     Answer: [Yes]

819     Justification: Large language models (LLMs) are used as important and non-standard  
820     components in both the AutoT2IS generation framework and the T2IS-Eval evaluation  
821     pipeline. In AutoT2IS, LLMs are employed during the semantic captioning stage to generate  
822     high-level semantic descriptions of generated images, which are then used to guide visual  
823     consistency alignment. In T2IS-Eval, LLMs are used to parse diverse user instructions and  
824     convert them into structured, multi-faceted evaluation criteria for consistency assessment.  
825     These uses of LLMs play a central role in ensuring instruction compliance and semantic  
826     alignment throughout the generation and evaluation process. The methodology and potential  
827     limitations of LLM usage are discussed in detail in the paper.

828     Guidelines:

- 829         • The answer NA means that the core method development in this research does not  
830         involve LLMs as any important, original, or non-standard components.

831  
832

- Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>)  
for what should or should not be described.