Generate Any Scene: Evaluating and Improving Text-to-Vision Generation with Scene Graph Programming

Anonymous CVPR submission

Abstract

001 Generative models like DALL-E and Sora have gained 002 attention by producing implausible images, such as "astronauts riding a horse in space." Despite the proliferation of 003 text-to-vision models that have inundated the internet with 004 synthetic visuals, from images to 3D assets, current bench-005 marks predominantly evaluate these models on real-world 006 007 scenes paired with captions. We introduce GENERATE ANY SCENE, a framework that systematically enumerates scene 008 graphs representing a vast array of visual scenes, spanning 009 010 realistic to imaginative compositions. GENERATE ANY 011 SCENE leverages 'scene graph programming,' a method for dynamically constructing scene graphs of varying complex-012 013 ity from a structured taxonomy of visual elements. This taxonomy includes numerous objects, attributes, and relations, 014 015 enabling the synthesis of an almost infinite variety of scene graphs. Using these structured representations, GENERATE 016 ANY SCENE translates each scene graph into a caption, en-017 abling scalable evaluation of text-to-vision models through 018 019 standard metrics. We conduct extensive evaluations across 020 multiple text-to-image, text-to-video, and text-to-3D mod-021 els, presenting key findings on model performance. We find that DiT-backbone text-to-image models align more closely 022 with input captions than UNet-backbone models. Text-to-023 video models struggle with balancing dynamics and consis-024 025 tency, while both text-to-video and text-to-3D models show 026 notable gaps in human preference alignment. Additionally, we demonstrate the effectiveness of GENERATE ANY 027 028 SCENE by conducting three practical applications leveraging captions generated by GENERATE ANY SCENE: (1) a 029 030 self-improving framework where models iteratively enhance 031 their performance using generated data, (2) a distillation 032 process to transfer specific strengths from proprietary models to open-source counterparts, and (3) improvements in 033 034 content moderation by identifying and generating challeng-035 ing synthetic data.

1. Introduction

Artist Marc Chagall said "Great art picks up where nature 037 ends." The charm of visual content generation lies in the 038 realm of imagination. Since their launch, Dall-E [4, 47] 039 and Sora [6] have promoted their products with implausible 040 generated images of "astronauts riding a horse in space" and 041 "cats playing chess". With the proliferation of text-to-vision 042 generation models, the internet is now flooded with gener-043 ated visual content-images, videos, and 3D assets-most 044 generated from user-provided captions [4, 6, 47]. While 045 there are numerous benchmarks designed for evaluating 046 these text-to-vision models, they are typically collections of 047 real-world visual content paired with captions [8, 31, 60]. 048 To quote Marc Chagall again, "If I create from heart, nearly 049 everything works; if from the head, almost nothing." There 050 is a need for evaluation benchmarks that go beyond real-051 world scenes and evaluate how well generative models can 052 represent the entire space of imaginary scenes. 053

036

072

073

074

Such a comprehensive evaluation requires that we first 054 define the space of the visual content. A long list of prior 055 work [23–25, 29, 40] has argued that scene graphs [29] are a 056 cognitively grounded [5] representation of the visual space. 057 A scene graph represents objects in a scene as individual 058 nodes in a graph. Each object is modified by attributes, 059 which describe its properties. For example, attributes can 060 describe the material, color, size, and location of the object 061 in the scene. Finally, relationships are edges that connect 062 the nodes. They define the spatial, functional, social, and 063 interactions between objects [37]. For example, in a liv-064 ing room scene, a "table" node might have attributes like 065 "wooden" or "rectangular" and be connected to a "lamp" 066 node through a relation: "on top of." This systematic scene 067 graph structure provides simple yet effective ways to define 068 and model the scene. Make it an ideal structure for GENER-069 ATE ANY SCENE to systematically define the diverse space 070 of the visual scenes. 071

We introduce GENERATE ANY SCENE, a system capable of efficiently enumerating the space of scene graphs representing a wide range of visual scenes, from realistic to

CVPR 2025 Submission . CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.



Figure 1. Overview of applications with GENERATE ANY SCENE captions: **Application 1**: (**Self-improving**): Iteratively enhances a model by generating images with GENERATE ANY SCENE captions, selecting the best, and fine-tuning, yielding a performance boost. **Application 2**: (**Distilling limitations**): Distills strengths from proprietary models, such as better compositionality and hard concept understanding, into open-source models. **Application 3**: (**Generated content detector**): Robustify AI-generated content detection by training on diverse synthetic data generated by GENERATE ANY SCENE's captions.

highly imaginative. GENERATE ANY SCENE is powered 075 076 by what we call *scene graph programming*, a programmatic approach for composing scene graphs of any complexity us-077 ing a rich taxonomy of visual elements, and for translating 078 each scene graph into a caption. With a space of syntheti-079 080 cally diverse captions, we use GENERATE ANY SCENE to prompt Text-to-Vision generation models and evaluate their 081 generations. Like any other representation, scene graphs 082 are also limited: they don't represent tertiary relationships 083 (e.g. "three people playing frisbee"). Nonetheless, they ac-084 085 count for a large space of possibilities. To systematically define and scalably explore the space of user captions, we 086 adopt the scene graph representation [29] to comprehen-087 sively evaluate and improve text-to-vision models. 088

We construct a rich taxonomy of visual concepts con-089 sisting of 28,787 objects, 1,494 attributes, 10,492 rela-090 tions, 2, 193 image/video/3D scene attributes from various 091 092 sources. Based on these assets, GENERATE ANY SCENE can programmatically synthesize an almost infinite num-093 ber of scene graphs of varying complexity [70]. Besides, 094 GENERATE ANY SCENE allows configurable scene graph 095 096 generation. For example, evaluators can specify the complexity level of the scene graph to be generated or pro-097 vide a seed scene graph to be expanded. Given an initial 098 scene graph, GENERATE ANY SCENE programmatically 099 translates it into a caption, which, when combined with 100 existing text-to-vision metrics, e.g., Clip Score [46] and 101 102 VOA Score [34], can be used to evaluate any text-to-vision model [53]. By automating these steps, our system ensures103both scalability and adaptability, providing researchers and104developers with diverse, richly detailed scene graphs and105corresponding captions tailored to their specific needs.106

With GENERATE ANY SCENE's programmatic genera-107 tion capability, we release a dataset featuring 10 million di-108 verse and compositional captions, each paired with a cor-109 responding scene graph. This extensive dataset spans a 110 wide range of visual scenarios, from realistic to highly 111 imaginative compositions, providing an invaluable resource 112 for researchers and practitioners in the Text-to-Vision gen-113 eration field. We also conduct extensive evaluations of 114 12 text-to-image, 9 text-to-video and 5 text-to-3D models 115 across a broad spectrum of visual scenes. We have several 116 crucial findings: (1) DiT-backbone models show superior 117 faithfulness and comprehensiveness to input captions than 118 UNet-backbone models, with human-alignment data train-119 ing helping to bridge some of these gaps. (2) Text-to-Video 120 generation face challenges in balancing dynamics and con-121 sistency. (3) All Text-to-Video and Text-to-3D models we 122 evaluate show negative ImageReward Score scores, high-123 lighting a substantial gap in human preference alignment. 124

Further, we demonstrate the effectiveness of GENER-
ATE ANY SCENE by conducting three practical applications125leveraging captions generated by GENERATE ANY SCENE126(Figure 1):128

Application 1: Self-improving. We show that our diverse129captions can facilitate a framework to iteratively improve130

Text-to-Vision generation models using their own genera-131 tions. Given a model, we generate multiple images, identify 132 133 the highest-scoring one, and use it as new fine-tuning data to improve the model itself. We fine-tune Stable Diffusion 134 135 v1-5 and achieve an average of 5% performance boost compared with original models, and this method is even better 136 than fine-tuning with the same amount of real images and 137 captions from the Conceptual Captions CC3M [8] over dif-138 139 ferent benchmarks.

Application 2: Distilling limitations. Using our evalua-140 141 tions, we identify limitations in open-sourced models that their proprietary counterparts excel at. Next, we distill these 142 specific capabilities from proprietary models. For exam-143 ple, DaLL-E 3 excels particularly in generating composite 144 images with multiple parts. We distill this capability into 145 146 Stable Diffusion v1-5, effectively bridging the gap between DaLL-E 3 and Stable Diffusion v1-5. 147

Application 3: Generated content detector. Content
moderation is a vital application, especially as *Text-to- Vision generation* models improve. We identify which kinds
of data content moderation models are bad at detecting, generate more of such content, and retrain the detectors. We
train a ViT-T with our generated data and boost its detection capabilities across benchmarks.

155 2. Generate Any Scene

We present our implementation of GENERATE ANY SCENE
system. (Figure 2) It programmatically synthesizes diverse
scene graphs in terms of both structure and content and
translates them into corresponding captions.

Scene graph. A scene graph is a structured representa-160 tion of a visual scene, where objects are represented as 161 nodes, their attributes (such as color and shape) are prop-162 erties of those nodes, and their relationships (such as spa-163 164 tial or semantic connections) are represented as edges. In recent years, scene graphs have played a crucial role in 165 visual understanding tasks, such as those found in Visual 166 Genome [29] and GQA [22] for visual question answering 167 (VQA). Their utility has expanded to various Text-to-Vision 168 generation tasks. For example, the DSG score [12] lever-169 170 ages MLMs to evaluate how well captions align with generated scenes by analyzing scene graphs. 171

Metadata Type	Number	Source
Objects	28,787	WordNet [39]
Attributes	1,494	Wikipedia [61], etc.
Relations	10,492	Robin [41]
Scene Attributes	2,193	Places365 [36], etc.

Table 1. Summary of the quantities and source of visual elements.

172 Taxonomy of visual elements. To construct a scene graph,
173 we use three main metadata types: objects, attributes, and

relations. We also have scene attributes that capture the 174 board aspect of the caption, such as art style, to create a 175 complete visual caption. The numbers and the source of our 176 metadata are illustrated in Table 1. Additionally, we build a 177 taxonomy that categorizes metadata into distinct levels and 178 types, enabling fine-grained analysis. This structure allows 179 for detailed assessments, such as evaluating model perfor-180 mance on "flower" as a general concept and on specific sub-181 categories like "daisy." More details in Appendix C. 182

183

197

198

199

200

201

202

203

204

205

206

207

208

2.1. Scene graph programming

Step 1: Scene graph structure enumeration and query. 184 Our system first generates and stores a variety of scene 185 graph structures based on a specified level of complexity, 186 defined by the total number of objects, relationships, and 187 attributes in each graph. The process begins by determin-188 ing the number of object nodes, and then by systematically 189 enumerating different combinations of relationships among 190 these objects and their associated attributes. Once all graph 191 structures meeting the complexity constraint are enumer-192 ated, they are stored in a database for later use. This enu-193 meration process is executed only once for each level of 194 complexity, allowing us to efficiently query the database for 195 suitable templates when needed. 196

Step 2: Populate the scene graph structure with metadata. Given a scene graph structure, the next step involves populating the graph with metadata. For each object node, attribute node, and relation edge, we sample the corresponding content from our metadata. This process is highly customizable: users can define the topics and types of metadata to be included (e.g., selecting only common metadata or specifying particular relationships between particular objects, among other options). By determining the scope of metadata sampling, we can precisely control the final content of the captions and easily extend the diversity and richness in the scene graphs by incorporating new datasets.

Step 3: Sampling scene attributes. In addition to scene 209 graphs that capture the visual content of the image, we also 210 include scene attributes that describe aspects such as the art 211 style, viewpoint, time span (for video), and 3D attributes 212 (for 3D content). These scene attributes are sampled di-213 rectly from our metadata, creating a list that provides con-214 textual details to enrich the description of the visual content. 215 Step 4: Translate scene graph to caption. We introduce 216 an algorithm that converts scene graphs and a list of scene 217 attributes into captions. The algorithm processes the scene 218 graph in topological order, transforming each object, its at-219 tributes, and relational edges into descriptive text. To main-220 tain coherence, it tracks each concept's occurrence, distin-221 guishing objects with identical names using terms like "the 222 first" or "the second." Objects that have been previously 223 referenced without new relations are skipped to avoid mis-224 referencing. This approach enhances caption clarity by pre-225



Figure 2. The generation pipeline of GENERATE ANY SCENE: Step 1: The system enumerates scene graph structures that contain objects, attributes, and relations based on complexity, and queries the corresponding scene graph structure that satisfies the needs. Step 2: It populates these structures with metadata, assigning specific content to each node. Scene graphs are completed in this step. Step 3: In addition to the scene graph, scene attributes—such as art style and camera settings—are sampled to provide contextual depth beyond the scene graph. Step 4: The GENERATE ANY SCENE system combines the scene graph and scene attributes, such as art style and camera settings, and then translates them into a coherent caption by organizing the elements into structured text.

venting repetition and maintaining a logical reference. 226

3. Evaluating Text-to-Vision generation models 227

3.1. Experiment Settings 228

- Details of experiment settings are in Appendix D. 229
- Models. We conduct experiments on 12 Text-to-image 230 models [1, 4, 10, 11, 15, 30, 32, 43, 44, 48], 9 Text-to-Video 231 models [9, 17, 27, 52, 55, 57, 63, 69, 71], and 5 Text-to-3D 232 models [33, 38, 45, 56, 59]. Text-to-image models are eval-233 uated at a resolution of 1024×1024 pixels. We standardize 234 235 the frame length to 16 across all Text-to-Video models for fair comparisons. For Text-to-3D, we generate videos by 236 237 rendering from 120 viewpoints.
- Metrics. Across all Text-to-Vision generation tasks, we use 238 Clip Score [7] (semantic similarity), VOA Score [34] (faith-239 240 fulness), TIFA Score [12, 20] (faithfulness), Pick Score [28] (human preference), and ImageReward Score [66] (human 241 242 preference) as general metrics, and for Text-to-Video generation, VBench [21] for fine-grained video analysis like 243 244 consistency and dynamics.
- Synthetic captions. We evaluate our Text-to-Image genera-245 tion and Text-to-Video generation models on 10K randomly 246 generated captions, with scene graph complexity ranging 247 from 3 to 12 and scene attributes from 0 to 5, using unre-248 stricted metadata. For Text-to-3D generation models, due 249 to their limitations in handling complex captions and time-250 251 intensive generation, we restrict scene graph complexity to

1-3, scene attributes to 0-2, and evaluate on 1K captions. 252

253

259

260

261

262

263

264

265

266

267

268

269

270

271

3.2. Overall results

We evaluate Text-to-Image generation, Text-to-Video gen-254 eration, and Text-to-3D generation models on GENERATE 255 ANY SCENE. Here, we only list key findings; more details 256 and raw results can be found in Appendix D. 257 258

Text-to-Image generation results. (Figure 3)

- 1. DiT-backbone models outperform UNet-backbone models on VOA Score and TIFA Score, indicating greater faithfulness and comprehensiveness to input captions.
- 2. Despite using a UNet architecture, Playground v2.5 achieves higher Pick Score and ImageReward Score scores than other open-source models. We attribute this to Playground v2.5's alignment with human preferences achieved during training.
- 3. The closed-source model DaLL-E 3 maintains a significant lead in VQA Score, TIFA Score, and ImageReward Score, demonstrating strong faithfulness and alignment with prompts across generated content.

Text-to-Video generation results. (Table 2,3)

- 1. Text-to-video models face challenges in balancing dy-272 namics and consistency (Table 3). This is especially evi-273 dent in Open-Sora 1.2, which achieves high consistency 274 but minimal dynamics, and Text2Video-Zero, which ex-275 cels in dynamics but suffers from frame inconsistency. 276
- 2. All models exhibit negative ImageReward Score (Ta-277 ble 2), suggesting a lack of human-preferred visual ap-278



Figure 3. Comparative evaluation of *Text-to-Image generation* models across different backbones (DiT and UNet) using multiple metrics: *TIFA Score*, *Pick Score*, *VQA Score*, and *ImageReward Score*.

Model clip score		pick score	image reward score	VQA score	TiFA score	
VideoCraft2 [9]	0.2398	[RGB]255, 255, 2000.1976	[RGB]255, 200, 200-0.4202	0.5018	[RGB]255, 200, 2000.2466	
AnimateLCM [55]	0.2450	[RGB]255, 200, 2000.1987	[RGB]255, 255, 200-0.5754	0.4816	0.2176	
AnimateDiff [17]	[RGB]255, 200, 2000.2610	0.1959	-0.7301	[RGB]255, 255, 2000.5255	0.2208	
Open-Sora 1.2 [71]	0.2259	0.1928	-0.6277	[RGB]255, 200, 2000.5519	[RGB]255, 255, 2000.2414	
FreeInit [63]	[RGB]255, 200, 2000.2579	0.1950	-0.9335	0.5123	0.2047	
ModelScope [57]	0.2041	0.1886	-1.9172	0.3840	0.1219	
Text2Video-Zero [27]	0.2539	0.1933	-1.2050	0.4753	0.1952	
CogVideoX-2B [69]	0.2038	0.1901	-1.2301	0.4585	0.1997	
ZeroScope [52]	0.2289	0.1933	-1.1599	0.4892	0.2388	

Table 2. Overall performance of *Text-to-Video generation* models over 10K GENERATE ANY SCENE captions. Red Cell is the highest score. Yellow Cell is the second highest score.

Model	subject consistency	background consistency	motion smoothness	dynamic degree
Open-Sora 1.2	[RGB]255, 200, 2000.9964	[RGB]255, 200, 2000.9907	[RGB]255, 200, 2000.9973	[RGB]200, 220, 2550.0044
Text2Video-Zero	[RGB]200, 220, 2550.8471	[RGB]200, 220, 2550.9030	[RGB]200, 220, 2550.8301	[RGB]255, 200, 2000.9999
VideoCraft2	0.9768	0.9688	0.9833	0.3556
AnimateDiff	0.9823	0.9733	0.9859	0.1406
FreeInit	0.9581	0.9571	0.9752	0.4440
ModelScope	0.9795	0.9831	0.9803	0.1281
AnimateLCM	0.9883	0.9802	0.9887	0.0612
CogVideoX-2B	0.9583	0.9602	0.9823	0.4980
ZeroScope	0.9814	0.9811	0.9919	0.1670

Table 3. Overall performance of *Text-to-Video generation* models over 10K GENERATE ANY SCENE captions with VBench metrics. Red Cell is the highest score. Blue Cell is the lowest score.

- peal in the generated content, even in cases where certainmodels demonstrate strong semantic alignment.
- 281 3. *VideoCrafter2* strikes a balance across key metrics, leading in human-preference alignment, faithfulness, consistency, and dynamic.
- **284** *Text-to-3D generation* results. (Table 4)

Model	clip score	pick score	vqa score	tifa score	image reward score
Latent-NeRF [38]	0.2115	0.1910	0.4767	0.2216	-1.5311
DreamFusion-sd [45]	0.1961	0.1906	0.4421	0.1657	-1.5582
Magic3D-sd [33]	0.1947	0.1903	0.4193	0.1537	-1.6327
SJC [56]	0.2191	0.1915	0.5015	0.2563	-1.4370
DreamFusion-IF [45]	0.1828	0.1857	0.3872	0.1416	-1.9353
Magic3D-IF [33]	0.1919	0.1866	0.4039	0.1537	-1.8465
ProlificDreamer [59]	0.2125	0.1940	0.5411	0.2704	-1.2774

Table 4. Overall performance of *Text-to-3D generation* models over 10K GENERATE ANY SCENE captions.

- 1. *ProlificDreamer* outperforms other models, particularly in *ImageReward Score*, *VQA Score* and *TIFA Score*. 286
- All models receive negative *ImageReward Score* scores, highlighting a significant gap between human preference and current *Text-to-3D generation* generation capabilities.
 280 280 280 280 280 280

4. Application 1: Self-Improving Models

291

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

In this section, we explore how GENERATE ANY SCENE 292 facilitates a self-improvement framework for model gen-293 eration capabilities. By programmatically generating scal-294 able compositional captions from scene graphs, GENERATE 295 ANY SCENE expands the textual and visual space, allow-296 ing for a diversity of synthetic images that extend beyond 297 real-world scenes. Our goal is to utilize these richly varied 298 synthetic images to further boost model performance. 299

Iterative self-improving framework. Inspired by Dream-Sync [53], we designed an iterative self-improving framework using GENERATE ANY SCENE with *Stable Diffusion* v1-5 as the baseline model. With *VQA Score*, which shows strong correlation with human evaluations on compositional images [34], we guide the model's improvement throughout the process.

Specifically, GENERATE ANY SCENE generates 3×10 K captions across three epochs. For each caption, *Stable Diffusion v1-5* generates 8 images, and the image with the highest *VQA Score* is selected. From each set of 10K optimal images, we then select the top 25% (2.5k image-caption pairs) as the training data for each epoch. In subsequent epochs, we use the fine-tuned model from the prior iteration to generate new images. We employ LoRA [19] for parameter-efficient fine-tuning. Additional details are available in Appendix E.

To evaluate the effectiveness of self-improvement using synthetic data generated by GENERATE ANY SCENE, 318 we conduct comparative experiments with the CC3M 319 dataset, which comprises high-quality and diverse realworld image-caption pairs [51]. We randomly sample $3 \times$ 321 10K captions from CC3M, applying the same top-score selection strategy for iterative fine-tuning of *Stable Diffusion* 323



Figure 4. **Results for Application 1: Self-Improving Models**. Average VQA score of *Stable Diffusion v1-5* fine-tuned on different data across 1K GENERATE ANY SCENE image/video evaluation set and GenAI-Bench image/video benchmark [31].

v1-5. Additionally, we include a baseline using randomsample fine-tuning strategy to validate the advantage of our
highest-scoring selection-based strategy.

Results. We evaluate our self-improving pipeline on *Text-to-Vision generation* benchmarks, including GenAI
Bench [31]. For the *Text-to-Video generation* task, we
use *Text2Video-Zero* as the baseline model, substituting its
backbone with the original *Stable Diffusion v1-5* and our
fine-tuned *Stable Diffusion v1-5* models.

333 Our results show that fine-tuning with GENERATE ANY SCENE-generated synthetic data consistently outperforms 334 CC3M-based fine-tuning across Text-to-Vision generation 335 tasks (Figure 4), achieving the highest gains with our 336 337 highest-scoring selection strategy. This highlights GENER-ATE ANY SCENE's scalability and compositional diversity, 338 339 enabling models to effectively capture complex scene structures. Additional results are in Appendix F. 340

Takeaway for application 1

Iterative self-improving *Text-to-Vision generation* models with compositional and diverse synthetic captions can surpass fine-tuning with real-world image-caption data.

Potential reason: The compositional, synthetic captions generated by GENERATE ANY SCENE exhibit greater diversity than real-world data.

341

5. Application 2: Distilling limitations

Although self-improving with GENERATE ANY SCENEgenerated data shows clear advantages over high-quality
real-world datasets, its efficiency remains inherently constrained by the limitations of the model's own generation
ability. To address this, we leverage the taxonomy and programmatic generation capabilities within GENERATE ANY

SCENE to identify specific strengths of proprietary mod-349els (*DaLL-E 3*), and to distill these capabilities into open-350source models. More details are in Appendix F.351

5.1. Fine-Grained Analysis of DaLL-E 3's Exceptional Performance

352

353

366

As shown in Figure 3. DaLL-E 3 achieves TIFA Score 1.5 354 to 2 times higher than those of other models. When we 355 compare TIFA Score across varying numbers of elements 356 (objects, relations, and attributes per caption) in Figure 6b, 357 DaLL-E 3 counterintuitively maintains consistent perfor-358 mance regardless of element count. The performance of 359 other models declines as the element count increases, which 360 aligns with expected compositional challenges. We suspect 361 these differences are primarily due to DaLL-E 3's advanced 362 capabilities in compositionality and understanding hard 363 concepts, which ensures high faithfulness across diverse 364 combinations of element types and counts. 365

5.2. Distilling compositionality from DaLL-E 3

Observations. We find that DaLL-E 3 tends to produce367straightforward combinations of multiple objects (Figure 5).368In contrast, open-source models like Stable Diffusion v1-5369often omit some objects from the captions, even though they
are capable of generating each object individually.371

This difference suggests that DaLL-E 3 may be trained 372 on datasets emphasizing multi-object presence without rig-373 orous attention to image layout or object interaction. Such 374 training likely underpins DaLL-E 3's stronger performance 375 on metrics like TIFA Score and VOA Score, prioritizing ob-376 ject inclusion over detailed compositional arrangement. 377 Finetuning. To encourage Stable Diffusion v1-5 to learn 378 compositional abilities similar to those of DaLL-E 3., we 379 select a set of 778 images generated by DaLL-E 3, each con-380 taining multiple objects, and utilize this dataset to fine-tune 381 Stable Diffusion v1-5. For the baseline, we randomly sam-382

CVPR 2025 Submission . CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.



Figure 5. Examples for Application 2: Distilling limitations. Examples of images generated by DaLL-E 3, the original Stable Diffusion v1-5, and the fine-tuned versions. The left four captions demonstrate fine-tuning with multi-object captions generated by GENERATE ANY SCENE for better compositionality, while the right two columns focus on understanding hard concepts.

383 pled an equivalent set of DaLL-E 3-generated images paired 384 with generated captions from GENERATE ANY SCENE. Results. To evaluate compositional improvements, we gen-385 erate 1K multi-object captions. Figure 6b shows a 10% 386 TIFA Score increase after fine-tuning, random fine-tuning 387

by an average of 3%. These results indicate enhanced compositional abilities in handling complex generation tasks.

388

389

391

390 We analyze images generated by Stable Diffusion v1-5 before and after fine-tuning on high-complexity image-392 caption pairs (Figure 5). It is surprising to see that, with only 1K LoRA fine-tuning steps, Stable Diffusion v1-5 ef-393 fectively learn DaLL-E 3 's capability to arrange and com-394 pose multiple objects within a single image,. This fine-395 396 tuning strategy notably enhances alignment between generated images and their given captions. 397

On a broader set of 10K GENERATE ANY SCENE-398 generated captions, the fine-tuned model consistently out-399 performed the randomly fine-tuned model (Figure 6a), con-400 firming the generalizability and superiority of targeted fine-401 tuning for improving model performance. 402

5.3. Learning hard concepts from DaLL-E 3 403

Observation. Figure 5 shows that is capable not only of 404 handling multi-object generation but also of understanding 405 406 and generating rare and hard concepts, such as a specific species of flower. We attribute this to its training with pro-407 prietary real-world data. 408

Finetuning. Using the taxonomy of GENERATE ANY 409 SCENE, we compute model performance on each con-410 411 ceptby averaging scores across captions containing that 412 concept.Accumulating results through the taxonomy, we

identify the 100 concepts where Stable Diffusion v1-5 413 shows the largest performance gap relative to DaLL-E 3. 414 For fine-tuning, we generate 778 captions incorporating 415 these concepts with others, using DaLL-E 3 to produce cor-416 responding images. As a baseline, we randomly select 778 417 **GENERATE ANY SCENE-generated captions for fine-tuning** 418 and compare these with the original Stable Diffusion v1-5 419 model. 420

Results. The results in Figure 6c show that our targeted fine-tuning led to improved model performance, reflected in higher average scores across captions and increased scores for each challenging concept.

Takeaway for application 2

Targeted fine-tuning can distill proprietary model strengths, effectively bridging gaps in compositionality and concept handling for open-source models.

Potential Reason: GENERATE ANY SCENE facilitates fine-grained analysis to identify specific performance gaps, enabling targeted data selection to distill limitations.

425

426

421

422

423

424

6. Application 3: Generated content detector

Advances in Text-to-Vision generation underscore the need 427 for effective content moderation [42]. Major challenges in-428 clude the lack of high-quality and diverse datasets and the 429 difficulty of generalizing detection across models Text-to-430 Vision generation [26, 58]. GENERATE ANY SCENE ad-431





(b) **Distilling compositionality from DaLL-E 3**: Model results on TIFA vs. total element numbers in captions in 1K multi-object GEN-ERATE ANY SCENE captions.



Figure 6. Results for Application 2: Distilling limitations. The left two figures show the results for Distilling compositionality from DALL-E 3, while the rightmost figure shows the results for Learning hard concepts from DALL-E 3.



Figure 7. Results for Application 3: Generated content detector. Comparison of detection performance across different data scales using D^3 alone versus the combined D^3 + GENERATE ANY SCENE training set in cross-model and cross-dataset scenarios.

dresses these issues by enabling scalable, programmatic
generation of compositional captions, increasing the diversity and volume of synthetic data. This approach enhances
existing datasets by compensating for their limited scopefrom realistic to imaginative-and variability.

To demonstrate GENERATE ANY SCENE's effectiveness 437 in training generated content detectors, we used the D^3 438 dataset [2] as a baseline. We sampled 5k captioned real and 439 SDv1.4-generated image pairs from D^3 and generated 5k 440 additional images with GENERATE ANY SCENE captions. 441 We trained a ViT-T [62] model with a single-layer linear 442 443 classifier, varying dataset sizes with N real and N synthetic images. For synthetic data, we compared N samples solely 444 from D^3 with a mixed set of N/2 from GENERATE ANY 445 SCENE and N/2 from D^3 , keeping the same training size. 446

We evaluate the detector's generalization on the GenImage [72] validation set and images generated using GENERATE ANY SCENE captions. Figure 7 demonstrates that
combining GENERATE ANY SCENE-generated images with
real-world captioned images consistently enhances detection performance, particularly across cross-model scenarios
and diverse visual scenes. More details are in Appendix G.

Takeaway for application 3

Compositional synthetic captions robustify generated content detectors.

Potential reason: GENERATE ANY SCENE can generate more diverse captions to complement real-world image-caption training data by enriching compositional variety and imaginative scope.

7. Conclusion

We present GENERATE ANY SCENE, a system leveraging 456 scene graph programming to generate diverse and composi-457 tional synthetic captions for Text-to-Vision generation tasks. 458 It extends beyond existing real-world caption datasets to in-459 clude imaginary scenes and even implausible scenarios. To 460 demonstrate the effectiveness of GENERATE ANY SCENE, 461 we explore three applications: (1) self-improvement by it-462 eratively optimizing models, (2) distillation of proprietary 463 model strengths into open-source models, and (3) robust 464 content moderation with diverse synthetic data. GENERATE 465 ANY SCENE highlights the importance of synthetic data in 466 evaluating and improving Text-to-Vision generation, and ad-467 dresses the need to systematically define and scalably pro-468 duce the space of visual scenes. 469

470 References

498

499

500

501

- 471 [1] DeepFloyd Lab at StabilityAI. DeepFloyd IF: a novel state-of-the-art open-source text-to-image model with a high degree of photorealism and language understanding. https: 474 //www.deepfloyd.ai/deepfloyd-if, 2023. Retrieved on 2023-11-08. 4, 16, 17
- 476 [2] Lorenzo Baraldi, Federico Cocchi, Marcella Cornia,
 477 Alessandro Nicolosi, and Rita Cucchiara. Contrasting deep478 fakes diffusion via contrastive learning and global-local sim479 ilarities. *arXiv preprint arXiv:2407.20337*, 2024. 8
- [3] Sean Bell, Paul Upchurch, Noah Snavely, and Kavita Bala.
 Material recognition in the wild with the materials in context database. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3479–3487, 2015. 16
- [4] James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jianfeng
 Wang, Linjie Li, Long Ouyang, Juntang Zhuang, Joyce
 Lee, Yufei Guo, et al. Improving image generation with
 better captions. *Computer Science. https://cdn. openai. com/papers/dall-e-3. pdf*, 2(3):8, 2023. 1, 4, 16, 17
- 490 [5] Irving Biederman. Recognition-by-components: a theory of human image understanding. *Psychological review*, 94(2):
 492 115, 1987. 1
- [6] Tim Brooks, Bill Peebles, Connor Holmes, Will DePue,
 Yufei Guo, Li Jing, David Schnurr, Joe Taylor, Troy Luhman, Eric Luhman, et al. Video generation models as world
 simulators. 2024. URL https://openai. com/research/videogeneration-models-as-world-simulators, 3, 2024. 1, 16
 - [7] Tuhin Chakrabarty, Kanishk Singh, Arkadiy Saakyan, and Smaranda Muresan. Learning to follow object-centric image editing instructions faithfully. *ArXiv*, abs/2310.19145, 2023.
 4
- 502 [8] Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu
 503 Soricut. Conceptual 12M: Pushing web-scale image-text
 504 pre-training to recognize long-tail visual concepts. In *CVPR*,
 505 2021. 1, 3, 18
- [9] Haoxin Chen, Yong Zhang, Xiaodong Cun, Menghan Xia,
 Xintao Wang, Chao Weng, and Ying Shan. Videocrafter2:
 Overcoming data limitations for high-quality video diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7310–
 7320, 2024. 4, 5, 16, 17
- [10] Junsong Chen, Jincheng Yu, Chongjian Ge, Lewei Yao, Enze
 Xie, Yue Wu, Zhongdao Wang, James Kwok, Ping Luo,
 Huchuan Lu, et al. Pixart-\alpha: Fast training of diffusion
 transformer for photorealistic text-to-image synthesis. *arXiv preprint arXiv:2310.00426*, 2023. 4, 16, 17
- [11] Junsong Chen, Chongjian Ge, Enze Xie, Yue Wu, Lewei
 Yao, Xiaozhe Ren, Zhongdao Wang, Ping Luo, Huchuan Lu,
 and Zhenguo Li. Pixart-\sigma: Weak-to-strong training of
 diffusion transformer for 4k text-to-image generation. *arXiv preprint arXiv:2403.04692*, 2024. 4, 16, 17
- [12] Jaemin Cho, Yushi Hu, Roopal Garg, Peter Anderson, Ranjay Krishna, Jason Baldridge, Mohit Bansal, Jordi Pont-Tuset, and Su Wang. Davidsonian scene graph: Improving reliability in fine-grained evaluation for text-to-image generation. *ArXiv*, abs/2310.18235, 2023. 3, 4

 [13] Mircea Cimpoi, Subhransu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi. Describing textures in the wild. In *Proceedings of the IEEE conference on computer* vision and pattern recognition, pages 3606–3613, 2014. 16

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

- [14] Colby Crawford. 1000 cameras dataset. https:// www.kaggle.com/datasets/crawford/1000cameras-dataset, 2018. Accessed: 2024-11-09. 16
- [15] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first International Conference on Machine Learning*, 2024. 4, 16, 17
- [16] Y.C. Guo, Y.T. Liu, R. Shao, C. Laforte, V. Voleti, G. Luo, C.H. Chen, Z.X. Zou, C. Wang, Y.P. Cao, and S.H. Zhang. threestudio: A unified framework for 3d content generation. https://github.com/threestudio-project/ threestudio, 2023. 16
- [17] Yuwei Guo, Ceyuan Yang, Anyi Rao, Zhengyang Liang, Yaohui Wang, Yu Qiao, Maneesh Agrawala, Dahua Lin, and Bo Dai. Animatediff: Animate your personalized textto-image diffusion models without specific tuning. *arXiv preprint arXiv:2307.04725*, 2023. 4, 5, 16, 17
- [18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceed-ings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 20
- [19] J. Edward Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *ArXiv*, abs/2106.09685, 2021. 5
- [20] Yushi Hu, Benlin Liu, Jungo Kasai, Yizhong Wang, Mari Ostendorf, Ranjay Krishna, and Noah A Smith. Tifa: Accurate and interpretable text-to-image faithfulness evaluation with question answering, 2023. 4
- [21] Ziqi Huang, Yinan He, Jiashuo Yu, Fan Zhang, Chenyang Si, Yuming Jiang, Yuanhan Zhang, Tianxing Wu, Qingyang Jin, Nattapol Chanpaisit, et al. Vbench: Comprehensive benchmark suite for video generative models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21807–21818, 2024. 4, 17
- [22] Drew A Hudson and Christopher D Manning. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6700–6709, 2019. 3
- [23] Jingwei Ji, Ranjay Krishna, Li Fei-Fei, and Juan Carlos Niebles. Action genome: Actions as compositions of spatiotemporal scene graphs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10236–10247, 2020. 1
- [24] Justin Johnson, Ranjay Krishna, Michael Stark, Li-Jia Li, David Shamma, Michael Bernstein, and Li Fei-Fei. Image retrieval using scene graphs. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3668–3678, 2015.
- [25] Justin Johnson, Agrim Gupta, and Li Fei-Fei. Image generation from scene graphs. In *Proceedings of the IEEE con-*

ference on computer vision and pattern recognition, pages 1219–1228, 2018. 1

585

586

- 587 [26] Achhardeep Kaur, Azadeh Noori Hoshyar, Vidya Saikrishna,
 588 Selena Firmin, and Feng Xia. Deepfake video detection:
 589 challenges and opportunities. *Artificial Intelligence Review*,
 590 57(6):1–47, 2024. 7
- [27] Levon Khachatryan, Andranik Movsisyan, Vahram Tade-vosyan, Roberto Henschel, Zhangyang Wang, Shant Navasardyan, and Humphrey Shi. Text2video-zero: Text-to-image diffusion models are zero-shot video generators. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15954–15964, 2023. 4, 5, 16, 17
- 597 [28] Yuval Kirstain, Adam Polyak, Uriel Singer, Shahbuland Ma598 tiana, Joe Penna, and Omer Levy. Pick-a-pic: An open
 599 dataset of user preferences for text-to-image generation,
 600 2023. 4
- [29] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson,
 Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. Visual genome:
 Connecting language and vision using crowdsourced dense
 image annotations. *International journal of computer vision*,
 123:32–73, 2017. 1, 2, 3
- [30] Black Forest Labs. Flux.1: Advanced text-to-image models,
 2024. Accessed: 2024-11-10. 4, 16, 17
- [31] Baiqi Li, Zhiqiu Lin, Deepak Pathak, Jiayao Emily Li,
 Xide Xia, Graham Neubig, Pengchuan Zhang, and Deva Ramanan. Genai-bench: A holistic benchmark for compositional text-to-visual generation. In *Synthetic Data for Com- puter Vision Workshop*@ *CVPR 2024*, 2024. 1, 6, 14
- [32] Daiqing Li, Aleks Kamko, Ehsan Akhgari, Ali Sabet, Linmiao Xu, and Suhail Doshi. Playground v2. 5: Three insights towards enhancing aesthetic quality in text-to-image generation. *arXiv preprint arXiv:2402.17245*, 2024. 4, 16, 17
- [33] Chen-Hsuan Lin, Jun Gao, Luming Tang, Towaki Takikawa,
 Xiaohui Zeng, Xun Huang, Karsten Kreis, Sanja Fidler,
 Ming-Yu Liu, and Tsung-Yi Lin. Magic3d: High-resolution
 text-to-3d content creation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,
 pages 300–309, 2023. 4, 5, 16, 17
- [34] Zhiqiu Lin, Deepak Pathak, Baiqi Li, Jiayao Li, Xide Xia,
 Graham Neubig, Pengchuan Zhang, and Deva Ramanan.
 Evaluating text-to-visual generation with image-to-text generation. *ArXiv*, abs/2404.01291, 2024. 2, 4, 5, 17
- [35] Jiacheng Liu, Wenya Wang, Dianzhuo Wang, Noah A.
 Smith, Yejin Choi, and Hannaneh Hajishirzi. Vera: A
 general-purpose plausibility estimation model for commonsense statements, 2023. 12
- 633 [36] Alejandro López-Cifuentes, Marcos Escudero-Vinolo, Jesús
 634 Bescós, and Álvaro García-Martín. Semantic-aware scene
 635 recognition. *Pattern Recognition*, 102:107256, 2020. 3, 16
- [37] Cewu Lu, Ranjay Krishna, Michael Bernstein, and Li FeiFei. Visual relationship detection with language priors. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Pro- ceedings, Part I 14*, pages 852–869. Springer, 2016. 1
- [38] Gal Metzer, Elad Richardson, Or Patashnik, Raja Giryes, and
 Daniel Cohen-Or. Latent-nerf for shape-guided generation

of 3d shapes and textures. In Proceedings of the IEEE/CVF643Conference on Computer Vision and Pattern Recognition,644pages 12663–12673, 2023. 4, 5, 16, 17645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

- [39] George A Miller. Wordnet: a lexical database for english. *Communications of the ACM*, 38(11):39–41, 1995. 3, 14
- [40] Julian Ost, Fahim Mannan, Nils Thuerey, Julian Knodt, and Felix Heide. Neural scene graphs for dynamic scenes. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 2856–2865, 2021. 1
- [41] Jae Sung Park, Zixian Ma, Linjie Li, Chenhao Zheng, Cheng-Yu Hsieh, Ximing Lu, Khyathi Chandu, Norimasa Kobori Quan Kong, Ali Farhadi, and Ranjay Krishna Yejin Choi. Robin: Dense scene graph generations at scale with improved visual reasoning. 2024. 3, 16
- [42] Gan Pei, Jiangning Zhang, Menghan Hu, Zhenyu Zhang, Chengjie Wang, Yunsheng Wu, Guangtao Zhai, Jian Yang, Chunhua Shen, and Dacheng Tao. Deepfake generation and detection: A benchmark and survey. arXiv preprint arXiv:2403.17881, 2024. 7
- [43] Pablo Pernias, Dominic Rampas, Mats L Richter, Christopher J Pal, and Marc Aubreville. Würstchen: An efficient architecture for large-scale text-to-image diffusion models. *arXiv preprint arXiv:2306.00637*, 2023. 4, 16, 17
- [44] 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. arXiv preprint arXiv:2307.01952, 2023. 4, 16, 17
- [45] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. arXiv preprint arXiv:2209.14988, 2022. 4, 5, 16, 17
- [46] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision, 2021. 2
- [47] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1 (2):3, 2022. 1, 16
- [48] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of* the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 10684–10695, 2022. 4, 16, 17
- [49] Chitwan 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. 16
- [50] Babak Saleh and Ahmed Elgammal. Large-scale classification of fine-art paintings: Learning the right metric on the right feature. *arXiv preprint arXiv:1505.00855*, 2015. 16
- [51] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of the 56th Annual Meeting of the Association for* 700

- Computational Linguistics (Volume 1: Long Papers), pages
 2556–2565, 2018. 5
- 703
 [52] Spencer Sterling. zeroscope_v2_576w, 2023. Accessed:

 704
 2024-11-10. 4, 5, 16, 17
- [53] Jiao Sun, Deqing Fu, Yushi Hu, Su Wang, Royi Rassin,
 Da-Cheng Juan, Dana Alon, Charles Herrmann, Sjoerd van
 Steenkiste, Ranjay Krishna, and Cyrus Rashtchian. Dreamsync: Aligning text-to-image generation with image understanding feedback. *ArXiv*, abs/2311.17946, 2023. 2, 5
- [54] Giuseppe Vecchio and Valentin Deschaintre. Matsynth:
 A modern pbr materials dataset. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22109–22118, 2024. 16
- [55] Fu-Yun Wang, Zhaoyang Huang, Xiaoyu Shi, Weikang Bian,
 Guanglu Song, Yu Liu, and Hongsheng Li. Animatelcm:
 Accelerating the animation of personalized diffusion models and adapters with decoupled consistency learning. *arXiv preprint arXiv:2402.00769*, 2024. 4, 5, 16, 17
- [56] Haochen Wang, Xiaodan Du, Jiahao Li, Raymond A Yeh,
 and Greg Shakhnarovich. Score jacobian chaining: Lifting
 pretrained 2d diffusion models for 3d generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision*and Pattern Recognition, pages 12619–12629, 2023. 4, 5,
 16, 17
- [57] Jiuniu Wang, Hangjie Yuan, Dayou Chen, Yingya Zhang,
 Xiang Wang, and Shiwei Zhang. Modelscope text-to-video
 technical report. *arXiv preprint arXiv:2308.06571*, 2023. 4,
 5, 16, 17
- [58] Tianyi Wang, Xin Liao, Kam Pui Chow, Xiaodong Lin, and
 Yinglong Wang. Deepfake detection: A comprehensive survey from the reliability perspective. ACM Computing Surveys, 2024. 7
- [59] Zhengyi Wang, Cheng Lu, Yikai Wang, Fan Bao, Chongxuan Li, Hang Su, and Jun Zhu. Prolificdreamer: High-fidelity and diverse text-to-3d generation with variational score distillation. *Advances in Neural Information Processing Systems*, 36, 2024. 4, 5, 16, 17
- [60] Zijie J. Wang, Evan Montoya, David Munechika, Haoyang
 Yang, Benjamin Hoover, and Duen Horng Chau. DiffusionDB: A large-scale prompt gallery dataset for text-toimage generative models. *arXiv:2210.14896 [cs]*, 2022. 1,
 16
- 743 [61] Wikipedia Contributors. Lists of colors. https://en.
 744 wikipedia.org/wiki/Lists_of_colors, 2024.
 745 Accessed: 2024-11-09. 3, 16
- [62] Kan Wu, Jinnian Zhang, Houwen Peng, Mengchen Liu, Bin Xiao, Jianlong Fu, and Lu Yuan. Tinyvit: Fast pretraining distillation for small vision transformers. In *European conference on computer vision*, pages 68–85. Springer, 2022. 8, 20
- [63] Tianxing Wu, Chenyang Si, Yuming Jiang, Ziqi Huang, and
 Ziwei Liu. Freeinit: Bridging initialization gap in video diffusion models. In *European Conference on Computer Vision*,
 pages 378–394. Springer, 2025. 4, 5, 16, 17
- [64] Xiaoshi Wu, Yiming Hao, Keqiang Sun, Yixiong Chen, Feng
 Zhu, Rui Zhao, and Hongsheng Li. Human preference score
 v2: A solid benchmark for evaluating human preferences of

text-to-image synthesis. *arXiv preprint arXiv:2306.09341*, 2023. 14

758

759

760

761

762

763

764

765

766

767

768

769

770

771

772

773

774

775

776

777

778

779

780

781

782

783

784

785

786

787

788

789

790

791

792

793

794

- [65] Xindi Wu, Dingli Yu, Yangsibo Huang, Olga Russakovsky, and Sanjeev Arora. Conceptmix: A compositional image generation benchmark with controllable difficulty. *ArXiv*, abs/2408.14339, 2024. 14
- [66] Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao Dong. Imagereward: Learning and evaluating human preferences for textto-image generation, 2023. 4
- [67] Zhe Xu, Dacheng Tao, Ya Zhang, Junjie Wu, and Ah Chung Tsoi. Architectural style classification using multinomial latent logistic regression. In *Computer Vision–ECCV 2014:* 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part I 13, pages 600–615. Springer, 2014. 16
- [68] Jia Xue, Hang Zhang, Kristin Dana, and Ko Nishino. Differential angular imaging for material recognition. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 764–773, 2017. 16
- [69] Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming Yang, Wenyi Hong, Xiaohan Zhang, Guanyu Feng, et al. Cogvideox: Text-to-video diffusion models with an expert transformer. *arXiv preprint arXiv:2408.06072*, 2024. 4, 5, 16, 17
- [70] Jieyu Zhang, Weikai Huang, Zixian Ma, Oscar Michel, Dong He, Tanmay Gupta, Wei-Chiu Ma, Ali Farhadi, Aniruddha Kembhavi, and Ranjay Krishna. Task me anything. In Advances in neural information processing systems, 2024. 2
- [71] Zangwei Zheng, Xiangyu Peng, Tianji Yang, Chenhui Shen, Shenggui Li, Hongxin Liu, Yukun Zhou, Tianyi Li, and Yang You. Open-sora: Democratizing efficient video production for all, 2024. 4, 5, 16, 17
- [72] Mingjian Zhu, Hanting Chen, Qiangyu Yan, Xudong Huang, Guanyu Lin, Wei Li, Zhijun Tu, Hailin Hu, Jie Hu, and Yunhe Wang. Genimage: A million-scale benchmark for detecting ai-generated image. Advances in Neural Information Processing Systems, 36, 2024. 8

Generate Any Scene: Evaluating and Improving Text-to-Vision Generation with Scene Graph Programming

Supplementary Material

796 A. More Analysis with GENERATE ANY 797 SCENE

With GENERATE ANY SCENE, we can generate infinitely
diverse and highly controllable prompts. Using GENERATE ANY SCENE, we conduct several analyses to provide
insights into the performance of today's *Text-to-Vision generation* models.

803 A.1. Performance analysis across caption properties

804 In this section, we delve into how model performance varies with respect to distinct properties of GENERATE ANY 805 SCENE captions. While GENERATE ANY SCENE is capable 806 of generating an extensive diversity of captions, these out-807 puts inherently differ in key characteristics that influence 808 809 model evaluation. Specifically, we examine three properties of the caption: Commonsense, Perplexity, and Scene 810 Graph Complexity (captured as the number of elements in 811 the captions). These properties are critical in understand-812 ing how different models perform across a spectrum of lin-813 guistic and semantic challenges presented by captions with 814 varying levels of coherence, plausibility, and compositional 815 816 richness.

817 Perplexity. (Figure 8) Perplexity is a metric used to mea818 sure a language model's unpredictability or uncertainty in
819 generating a text sequence. A higher perplexity value indi820 cates that the sentences are less coherent or less likely to be
821 generated by the model.

As shown in Figure 8, From left to right, when perplexity increases, indicating that the sentences become less reasonable and less typical of those generated by a language model, we observe no clear or consistent trends across all models and metrics. This suggests that the relationship between perplexity and model performance varies depending on the specific model and evaluation metric.

829 Commonsense. (Figure 9) Commonsense is an inherent
830 property of text. We utilize the Vera Score [35], a metric
831 generated by a fine-tuned LLM to evaluate the text's commonsense level.

As shown in Figure 9, from left to right, as the Vera Score increases—indicating that the captions exhibit greater commonsense reasoning—we observe a general improvement in performance across all metrics and models, except for *Clip Score*. This trend underscores the correlation between commonsense-rich captions and enhanced model 838 performance. 839

840

841

842

843

844

857

Element Numbers (Complexity of Scene Graph). (Figure 10) Finally, we evaluate model performance across total element numbers in the captions, which represent the complexity of scene graphs (objects + attributes + relations).

From left to right, the complexity of scene graphs be-845 comes higher, reflecting more compositional and intricate 846 captions. Across most metrics and models, we observe 847 a noticeable performance decline as the scene graphs be-848 come more complex. However, an interesting exception is 849 observed in the performance of DaLL-E 3. Unlike other 850 models, DaLL-E 3 performs exceptionally well on VQA 851 Score and TIFA Score, particularly on VQA Score, where 852 it even shows a slight improvement as caption complexity 853 increases. This suggests that DaLL-E 3 may have a unique 854 capacity to handle complex and compositional captions ef-855 fectively. 856

A.2. Analysis on different metrics

Compared with most LLM and VLM benchmarks that use858multiple-choice questions and accuracy as metrics. There is859no universal metric in evaluating *Text-to-Vision generation*860models. Researchers commonly used model-based metrics861like *Clip Score*, *VQA Score*, etc. Each of these metrics862is created and fine-tuned for different purposes with bias.863Therefore, we also analysis on different metrics.864

Clip Score isn't a universal metric. Clip Score is one 865 of the most widely used metrics in Text-to-Vision gener-866 ation for evaluating the alignment between visual content 867 and text. However, our analysis reveals that Clip Score is 868 not a perfect metric and displays some unusual trends. For 869 instance, as shown in Figures 8, 9, and 10, we compute 870 the perplexity across 10k prompts used in our study, where 871 higher perplexity indicates more unpredictable or disorga-872 nized text. Interestingly, unlike other metrics, Clip Score 873 decreases as perplexity lowers, suggesting that Clip Score 874 tends to favor more disorganized text. This behavior is 875 counterintuitive and highlights the potential limitations of 876 using Clip Score as a robust alignment metric. 877

Limitations of human preference-based metrics. We use two metrics fine-tuned using human preference data: 879



Figure 8. Average performance of models across different percentiles of perplexity of captions, evaluated on various metrics. From left to right, the perplexity decreases, indicating captions that are progressively more reasonable and easier for the LLM to generate.

880 Pick Score and ImageReward Score. However, we found 881 that these metrics exhibit a strong bias toward the data on 882 which they were fine-tuned. For instance, as shown in Table 5, Pick Score assigns similar scores across all models, 883 884 failing to provide significant differentiation or meaningful insights into model performance. In contrast, ImageReward 885 886 Score demonstrates clearer preferences, favoring models such as DaLL-E 3 and Playground v2.5, which incorporated 887 human-alignment techniques during their training. How-888 ever, this metric shows a significant drawback: it assigns 889 disproportionately large negative scores to models like Sta-890 *ble Diffusion v2-1*, indicating a potential over-sensitivity to 891 alignment mismatches. Such behavior highlights the limi-892 tations of these metrics in providing fair and unbiased eval-893 uations across diverse model architectures. 894

VQA Score and *TIFA Score* are relative reliable metrics.
Among the evaluated metrics, *VQA Score* and *TIFA Score*stand out by assessing model performance on VQA tasks,
rather than relying solely on subjective human preferences.
This approach enhances the interpretability of the evaluation process. Additionally, we observed that the results
from *VQA Score* and *TIFA Score* show a stronger corre-

lation with other established benchmarks. Based on these902advantages, we recommend prioritizing these two metrics903for evaluation. However, it is important to note that their904effectiveness is constrained by the limitations of the VQA905models utilized in the evaluation.906

907

A.3. Fairness analysis

We evaluate fairness by examining the model's performance908across different genders and races.Specifically, we cal-culate the average performance for each node and its as-910sociated child nodes within the taxonomy tree constructed911for objects.For example, the node "females" includeschild nodes such as "waitresses," and their combined per-913formance is considered in the analysis.914

Gender.In gender, we observe a notable performance gap915between females and males, as could be seen from Fig-
ure 11, Models are better at generating male concepts.916

Race.There are also performance gaps in different races.918From Figure 12, we found that "white (person)" and "black919(person)" perform better than "asian (person)", "Indian920(amerindian)", and "Latin American".921



Figure 9. Average performance of models across different percentiles of Vera Score for captions, evaluated on various metrics. From left to right, the Vera Score decreases, indicating captions that exhibit less commonsense reasoning and are more likely to describe implausible scenes.

B. Correlation of GENERATE ANY SCENE with other *Text-to-Vision generation* benchmarks

The GENERATE ANY SCENE benchmark uniquely relies
entirely on synthetic captions to evaluate models. To assess
the transferability of these synthetic captions, we analyzed
the consistency in model rankings across different benchmarks [31, 64, 65]. Specifically, we identified the overlap
of models evaluated by two benchmarks and computed the
Spearman correlation coefficient between their rankings.

As shown in the figure 13, GENERATE ANY SCENE 932 demonstrates a strong correlation with other benchmarks, 933 934 such as Conceptmix [65] and GenAI Bench [31], indicating the robustness and reliability of GENERATE ANY SCENE's 935 synthetic caption-based evaluations. This suggests that the 936 synthetic captions generated by GENERATE ANY SCENE 937 can effectively reflect model performance trends, aligning 938 closely with those observed in benchmarks using real-world 939 940 captions or alternative evaluation methods.

C. Details of Taxonomy of Visual Concepts

941

To construct a scene graph, we utilize three primary types942of metadata: objects, attributes, and relations, which rep-943resent the structure of a visual scene. Additionally, scene944attributes—which include factors like image style, perspec-945tive, and video time span—capture broader aspects of the946visual content. Together, the scene graph and scene at-947tributes form a comprehensive representation of the scene.948

Our metadata is further organized using a well-defined949taxonomy, enhancing the ability to generate controllable950prompts. This hierarchical taxonomy not only facilitates951the creation of diverse scene graphs, but also enables fine-
grained and systematic model evaluation.953

Objects. To enhance the comprehensiveness and taxon-
omy of object data, we leverage noun synsets and the struc-
ture of WordNet [39]. In WordNet, a *physical object* is de-
fined as "a tangible and visible entity; an entity that can
cast a shadow." Following this definition, we designate the
physical object as the root node, constructing a hierarchical
959
tree with all 28,787 hyponyms under this category as the set954
955



Figure 10. Average performance of models across different numbers of elements (objects + attributes + relations) in the scene graph (complexity of the scene graph) of the captions, evaluated on various metrics. From left to right, as the number of elements (complexity) increases, the scene graphs become more complicated and compositional.



Figure 11. Average performance scores of all models across different genders were evaluated using various metrics.

0.60 0.55 0.50 0.45 0.40 0.35 0.30 Black (0.25 Asian (person 0.20 Score Type

Performance Across Different Races

Figure 12. Average performance scores of all models across different races evaluated using various metrics.

of objects in our model. 961

962 Following WordNet's hypernym-hyponym relationships, we establish a tree structure, linking each object to its pri-963 964 mary parent node based on its first-listed hypernym. For objects with multiple hypernyms, we retain only the primary 965 966 parent to simplify the hierarchy. Furthermore, to reduce am-

biguity, if multiple senses of a term share the same parent, 967 we exclude that term itself and reassign its children to the original parent node. This approach yields a well-defined and disambiguated taxonomy.



Figure 13. Correlation of GENERATE ANY SCENE with other popular *Text-to-Vision generation* benchmarks.

971 Attributes. The attributes of a scene graph represent 972 properties or characteristics associated with each object. 973 We classify these attributes into *nine* primary categories. For color, we aggregate 677 unique entries sourced from 974 975 Wikipedia [61]. The *material* category comprises 76 types, referenced from several public datasets [3, 54, 68]. The 976 977 *texture* category includes 42 kinds from the Describable Textures Dataset [13], while the *architectural style* encom-978 passes 25 distinct styles [67]. Additionally, we collect 85 979 states, 41 shapes, and 24 sizes. For human descriptors, we 980 compile 59 terms across subcategories, including body type 981 and height. Finally, we collect 465 common adjectives cov-982 983 ering general characteristics of objects to enhance the de-984 scriptive richness of our scene graphs.

Relationships. We leverage the Robin dataset [41] as the 985 986 foundation for relationship metadata, encompassing six key categories: spatial, functional, interactional, social, emo-987 tional, and symbolic. With 10,492 relationships, the dataset 988 989 provides a comprehensive and systematic repository that supports modeling diverse and complex object interactions. 990 991 Its extensive coverage captures both tangible and abstract connections, forming a robust framework for accurate scene 992 993 graph representation.

994 Scene Attributes. In *Text-to-Vision generation* tasks,
995 people mainly focus on creating realistic images and art
996 from a text description [4, 47, 49]. For artistic styles,
997 we define scene attributes using 76 renowned *artists*, 41
998 genres, and 126 painting styles from WikiArt [50], along
999 with 29 common painting techniques. For realistic im1000 agery, we construct camera settings attributes across 6 cat-

egories: camera models, focal lengths, perspectives, aper-1001 tures, depths of field, and shot scales. The camera models 1002 are sourced from the 1000 Cameras Dataset [14], while the 1003 remaining categories are constructed based on photography 1004 knowledge and common prompts in Text-to-Vision genera-1005 *tion* tasks [6, 60]. To control scene settings, we categorize 1006 location, weather and lighting attributes, using 430 diverse 1007 locations from Places365 [36], alongside 76 weathers and 1008 57 lighting conditions. For video generation, we introduce 1009 attributes that describe dynamic elements. These include 1010 12 types of camera rig, 30 distinct camera movements, 15 1011 video editing styles, and 27 temporal spans. The compre-1012 hensive scene attributes that we construct allow for the de-1013 tailed and programmatic Text-to-Vision generation genera-1014 tion. 1015

D. Details of Overall Performance (Section 3) 1016

1017

1029

1030

1031

1032

1033

1034

D.1. Detailed experiment settings

- For Text-to-Image generation, we select a range of open-1018 source models, including those utilizing UNet back-1019 bones, such as DeepFloyd IF [1], Stable Diffusion v2-1020 1 [48], SDXL [44], Playground v2.5 [32], and Wuer-1021 stchen v2 [43], as well as models with DiT backbones, 1022 including Stable Diffusion 3 Medium [15], PixArt- α [10], 1023 PixArt- Σ [11], FLUX.1-schnell [30], FLUX.1-dev [30], 1024 and FLUX 1. Closed-source models, such as DaLL-E 1025 3 [4] and FLUX1.1 PRO [30], are also assessed to ensure 1026 a comprehensive comparison. All models are evaluated at 1027 a resolution of 1024×1024 pixels. 1028
- For *Text-to-Video generation*, we select nine open-source models: *ModelScope* [57], *ZeroScope* [52], *Text2Video-Zero* [27], *CogVideoX-2B* [69], *VideoCrafter2* [9], *AnimateLCM* [55], *AnimateDiff* [17], *FreeInit* [63], and *Open-Sora 1.2* [71]. We standardize the frame length to 16 across all video models for fair comparisons.
- For Text-to-3D generation, we evaluate five recently 1035 SJC [56], DreamFusion [45], proposed models: 1036 Magic3D [33], Latent-NeRF [38], and Prolific-1037 Dreamer [59]. We employ the implementation and 1038 configurations provided by ThreeStudio [16] and gen-1039 erate videos by rendering from 120 viewpoints. То 1040 accelerate inference, we omit the refinement stage. 1041 For Magic3D and DreamFusion, we respectively use 1042 DeepFloyd IF and Stable Diffusion v2-1 as their 2D 1043 backbones. 1044

Metrics.Across all Text-to-Image generation, Text-to-1045Video generation, and Text-to-3D generation, we employ1046five widely used Text-to-Vision generation metrics to com-1047prehensively assess model performance:1048

 Clip Score: Assesses semantic similarity between images and text.
 1049
 1050

- *VQA Score* and *TIFA Score*: Evaluate faithfulness by generating question-answer pairs and measuring answer accuracy from images.
 - *Pick Score* and *ImageReward Score*: Capture human preference tendencies.

We also use metrics in VBench [21] to evaluate *Text-to-Video generation* models on fine-grained dimensions, such as consistency and dynamics, providing detailed insights into video performance.

For *Text-to-Video generation* and *Text-to-3D generation* tasks:

- We calculate *Clip Score*, *Pick Score*, and *ImageReward Score* on each frame, then average these scores across all
 frames to obtain an overall video score.
- For VQA Score and TIFA Score, we handle Text-to-Video
 generation and Text-to-3D generation tasks differently:
- In *Text-to-Video generation* tasks, we uniformly sample
 four frames from the 16-frame sequence and arrange
 them in a 2 × 2 grid image.
- For *Text-to-3D generation* tasks, we render images at
 45-degree intervals from nine different viewpoints and
 arrange them in a 3 × 3 grid.
- 1073 This sampling approach optimizes inference speed with-1074 out affecting score accuracy [34].

1075 D.2. Detailed overall results

1054 1055

1056

1057

1058

1059

1060

1061

1076 We evaluate *Text-to-Image generation*, *Text-to-Video generation*, and *Text-to-3D generation* models on GENERATE
1078 ANY SCENE. The detailed results of each model and each
1079 metric are shown in Tabs. 5 to 8

Model	clip score	pick score	vqa score	tifa score	image reward score
Playground v2.5 [32]	0.2581	0.2132	0.5734	0.2569	0.2919
Stable Diffusion v2-1 [48]	0.2453	0.1988	0.5282	0.2310	-0.9760
SDXL [44]	0.2614	0.2046	0.5328	0.2361	-0.3463
Wuerstchen v2 [43]	0.2448	0.2022	0.5352	0.2239	-0.3339
DeepFloyd IF XL [1]	0.2396	0.1935	0.5397	0.2171	-0.8687
Stable Diffusion 3 Medium [15]	0.2527	0.2027	0.5579	0.2693	-0.0557
PixArt-a [10]	0.2363	0.2050	0.6049	0.2593	0.1149
PixArt-Σ [11]	0.2390	0.2068	0.6109	0.2683	0.0425
FLUX.1-dev [30]	0.2341	0.2060	0.5561	0.2295	0.1588
FLUX.1-schnell [30]	0.2542	0.2047	0.6132	0.2833	0.1251
FLUX1.1 PRO [30]	0.2315	0.2065	0.5744	0.2454	-0.0361
Dalle-3 [4]	0.2518	0.2006	0.6871	0.4249	0.3464

Table 5. Overall performance of *Text-to-Image generation* models over 10K GENERATE ANY SCENE prompts.

Model	clip score	pick score	image reward score	VQA score	TiFA score
VideoCraft2 [9]	0.2398	0.1976	-0.4202	0.5018	0.2466
AnimateDiff [17]	0.2610	0.1959	-0.7301	0.5255	0.2208
Open-Sora 1.2 [71]	0.2259	0.1928	-0.6277	0.5519	0.2414
FreeInit [63]	0.2579	0.1950	-0.9335	0.5123	0.2047
ModelScope [57]	0.2041	0.1886	-1.9172	0.3840	0.1219
Text2Video-Zero [27]	0.2539	0.1933	-1.2050	0.4753	0.1952
AnimateLCM [55]	0.2450	0.1987	-0.5754	0.4816	0.2176
CogVideoX-2B [69]	0.2038	0.1901	-1.2301	0.4585	0.1997
ZeroScope [52]	0.2289	0.1933	-1.1599	0.4892	0.2388

Table 6. Overall performance of *Text-to-Video generation* models over 10k GENERATE ANY SCENE prompts.

Model	subject consistency	background consistency	motion smoothness	dynamic degree	aesthetic quality	imaging quality
VideoCraft2	0.9768	0.9688	0.9833	0.3556	0.5515	0.6974
AnimateDiff	0.9823	0.9733	0.9859	0.1406	0.5427	0.5830
Open-Sora 1.2	0.9964	0.9907	0.9973	0.0044	0.5235	0.6648
FreeInit	0.9581	0.9571	0.9752	0.4440	0.5200	0.5456
ModelScope	0.9795	0.9831	0.9803	0.1281	0.3993	0.6494
Text2Video-Zero	0.8471	0.9030	0.8301	0.9999	0.4889	0.7018
AnimateLCM	0.9883	0.9802	0.9887	0.0612	0.6323	0.6977
CogVideoX-2B	0.9583	0.9602	0.9823	0.4980	0.4607	0.6098
ZeroScope	0.9814	0.9811	0.9919	0.1670	0.4582	0.6782

Table 7. Overall performance of *Text-to-Image generation* models over 10k GENERATE ANY SCENE prompts with VBench metrics.

Model	clip score	pick score	vqa score	tifa score	image reward score
ProlificDreamer [59]	0.2125	0.1940	0.5411	0.2704	-1.2774
Latent-NeRF [38]	0.2115	0.1910	0.4767	0.2216	-1.5311
DreamFusion-sd [45]	0.1961	0.1906	0.4421	0.1657	-1.5582
Magic3D-sd [33]	0.1947	0.1903	0.4193	0.1537	-1.6327
SJC [56]	0.2191	0.1915	0.5015	0.2563	-1.4370
DreamFusion-IF [45]	0.1828	0.1857	0.3872	0.1416	-1.9353
Magic3D-IF [33]	0.1919	0.1866	0.4039	0.1537	-1.8465

Table 8. Overall performance of *Text-to-3D generation* models over 10k GENERATE ANY SCENE prompts.

D.3. Case study: Pairwise fine-grained model comparison 1080

Evaluating models using a single numerical average score 1082 can be limiting, as different training data often lead models 1083 to excel in generating different types of concepts. By lever-1084 aging the taxonomy we developed for GENERATE ANY 1085 SCENE, we can systematically organize these concepts and 1086 evaluate each model's performance on specific concepts 1087 over the taxonomy. This approach enables a more de-1088 tailed comparison of how well models perform on individ-1089 ual concepts rather than relying solely on an overall aver-1090 age score. Our analysis revealed that, while the models 1091 may achieve similar average performance, their strengths 1092 and weaknesses vary significantly across different concepts. 1093 Here we present a pairwise comparison of models across 1094 different metrics. 1095



Figure 14. *Stable Diffusion v2-1* vs. *Stable Diffusion 3 Medium* on average *VQA Score* in fine-grained categories.



Figure 15. *PixArt*- Σ vs. *Stable Diffusion 3 Medium* on average *VQA Score* in fine-grained categories.



Figure 16. *FLUX.1-schnell* vs. *Stable Diffusion 3 Medium* on average *VQA Score* in fine-grained categories.



Figure 17. PixArt- Σ vs. FLUX.1-schnell on average VQA Score in fine-grained categories.

E. Details of Application 1: Self-ImprovingModels (Section 4)

1098 E.1. Experiment details

1099 E.1.1 Captions Preparation

To evaluate the effectiveness of our iterative self-improving *Text-to-Vision generation* model, we generated three distinct sets of 10k captions using GENERATE ANY SCENE,
covering a sample complexity range from 3 to 12. These
captions were programmatically created to reflect a spectrum of structured scene graph compositions, designed to

challenge and enrich the model's learning capabilities.

For comparative analysis, we leveraged the Concep-
tual Captions (CC3M) [8] dataset, a large-scale benchmark1107tual Captions (CC3M) [8] dataset, a large-scale benchmark1108containing approximately 3.3 million image-caption pairs1109sourced from web alt-text descriptions. CC3M is renowned1110for its diverse visual content and natural language expressions, encompassing a wide range of styles, contexts, and1112semantic nuances.1113

1106

1135

1136

1137

1138

1139

1147

To ensure fair comparison, we randomly sampled three 1114 subsets of 10k captions from the CC3M dataset, matching 1115 the GENERATE ANY SCENE-generated caption sets in size. 1116 This approach standardizes data volume while enabling di-1117 rect performance evaluation. The diversity and semantic 1118 richness of the CC3M captions serve as a robust benchmark 1119 to assess whether GENERATE ANY SCENE-generated cap-1120 tions can match or exceed the descriptive quality of real-1121 world data across varied visual contexts. 1122

E.1.2 Dataset Construction and Selection Strategies 1123

For the captions generated by GENERATE ANY SCENE, we 1124 employed a top-scoring selection strategy to construct the 1125 fine-tuning training dataset, using a random selection strat-1126 egy as a baseline for comparison. Specifically, for each 1127 prompt, the model generated eight images. Under the top-1128 scoring strategy, we evaluated the generated images using 1129 the VQA score and selected the highest-scoring image as 1130 the best representation of the prompt. This process yielded 1131 10k top-ranked images per iteration, from which the top 1132 25% (approximately 2.5k images) with the highest VQA 1133 scores were selected to form the fine-tuning dataset. 1134

In the random selection strategy, one image was randomly chosen from the eight generated per prompt, and 25% of these 10k randomly selected images were sampled to create the fine-tuning dataset, maintaining parity in data size.

For the CC3M dataset, each prompt was uniquely paired1140with a real image. From the 10k real image-caption pairs1141sampled from CC3M, the top 25% with the highest VQA1142scores were selected as the fine-tuning training dataset. This1143ensured consistency in data size and selection criteria across1144all methods, facilitating a rigorous and equitable comparison of fine-tuning strategies.1146

E.1.3 Fine-tuning details

We fine-tuned the *Stable Diffusion v1-5* using the LoRA 1148 technique. The training was conducted with a resolution of 512 × 512 for input images and a batch size of 8. Gradients 1149 were accumulated over two steps. The optimization process 1151 utilized the AdamW optimizer with $\beta_1 = 0.9$, $\beta_2 = 0.999$, 1152 an ϵ value of 1×10^{-8} , and a weight decay of 10^{-2} . The 1153 learning rate was set to 1×10^{-4} and followed a cosine 1154

scheduler for smooth decay during training. To ensure sta-1155 bility, a gradient clipping threshold of 1.0 was applied. The 1156 fine-tuning process was executed for one epoch, with a max-1157 imum of 2500 training steps. For the LoRA-specific config-1158 1159 urations, we set the rank of the low-rank adaptation layers and the scaling factor α to be 128. 1160

After completing fine-tuning for each epoch, we set the 1161 LoRA weight to 0.75 and integrate it into Stable Diffusion 1162 1163 v1-5 to guide image generation and selection for the next subset. For the CC3M dataset, images from the subsequent 1164 1165 subset are directly selected.

In the following epoch, the fine-tuned LoRA parame-1166 ters from the previous epoch are loaded and used to resume 1167 training on the current subset, ensuring continuity and lever-1168 aging the incremental improvements from prior iterations. 1169

E.2. More results of fine-tuning models 1170

1171 Aside from our own test set and GenAI benchmark, we also evaluated our fine-tuned Text-to-Image generation models 1172 on the Tifa Bench (Figure 18), where we observed the same 1173 trend: models fine-tuned with our prompts consistently out-1174 performed the original Stable Diffusion v1-5 and CC3M 1175 1176 fine-tuned models.



Figure 18. Results for Application 1: Self-Improving Models. Average TIFA score of Stable Diffusion v1-5 fine-tuned with different data over TIFA Bench.

F. Details of Application 2: Distilling limita-1177 tions (Section 5) 1178

F.1. Collecting hard concepts 1179

We selected 81 challenging object concepts where Stable 1180 Diffusion v1-5 and DaLL-E 3 exhibit the largest gap in VQA 1181 Score. To determine the score for each concept, we calcu-1182 lated the average VQA score of the captions containing that 1183 specific concept. The full list of hard concepts is shown 1184 below: 1185

- 1. cloverleaf 1186
- 2. aerie (habitation) 1187
- 3. admixture 1188
- 4. webbing (web) 1189
- 1190 5. platter

~		
6.	voussoir	1191
7.	hearthstone	1192
8.	puttee	1193
9.	biretta	1194
10.	yarmulke	1195
11.	surplice	1196
12.	overcoat	1197
13.	needlepoint	1198
14.	headshot	1199
15.	photomicrograph	1200
16.	lavaliere	1201
17.	crepe	1202
18.	tureen	1203
19.	bale	1204
20.	jetliner	1205
21.	square-rigger	1206
22.	supertanker	1207
23.	pocketcomb	1208
24.	filament (wire)	1209
25.	inverter	1210
26.	denture	1211
27.	lidar	1212
28	volumeter	1213
29	colonoscope	1214
30	synchroevelotron	1215
31	miller (shaper)	1216
32	alternator	1213
33	dicer	1218
34	trundle	1210
35	naddle (blade)	1213
36	harmonica	1220
37	niccolo	1221
38	handrest	1222
30.	rundle	1223
39. 40	blowtorch	1224
40.	vollevhall	1223
41.	tile (men)	1220
42.	chuttlessel	1227
45.		1228
44.	Jigsaw	1229
45.	roaster (pan)	1230
46.	maze	1231
4/.	belt (ammunition)	1232
48.		1233
49.	drawer (container)	1234
50.	tenter	1235
51.	pinnacie (steepie)	1236
52.	pegboard	1237
53.	atterdeck	1238
54.	scaffold	1239
55.	catheter	1240
56.	broomcorn	1241
57.	spearmint	1242
58.	okra (herb)	1243

1244	59.	goatsfoot
1245	60.	peperomia
1246	61.	ammobium
1247	62.	gazania
1248	63.	echinocactus
1249	64.	birthwort
1250	65.	love-in-a-mist (passionflower)
1251	66.	ragwort
1252	67.	spicebush (allspice)
1253	68.	leadplant
1254	69.	barberry
1255	70.	hamelia
1256	71.	jimsonweed
1257	72.	undershrub
1258	73.	dogwood
1259	74.	butternut (walnut)
1260	75.	bayberry (tree)
1261	76.	lodestar
1262	77.	tapa (bark)
1263	78.	epicalyx
1264	79.	blackberry (berry)
1265	80.	stub
1266	81.	shag (tangle)

1267 F.2. Experiment details

We conducted targeted fine-tuning experiments on Stable 1268 Diffusion v1-5 to evaluate GENERATE ANY SCENE's ef-1269 1270 fectiveness in distilling model compositionality and learn-1271 ing hard concepts. For each task, we selected a dataset of 778 GENERATE ANY SCENE captions paired with im-1272 1273 ages generated by DaLL-E 3. For compositionality, we selected multi-object captions from the existing dataset of 10k 1274 1275 GENERATE ANY SCENE captions and paired them with the 1276 corresponding images generated by DaLL-E 3. To address hard concept learning, we first used Stable Diffusion v1-5 to 1277 generate images based on the 10k GENERATE ANY SCENE 1278 captions and identified the hard concepts with the lowest 1279 1280 VQA scores. These concepts were then used to create a sub-1281 set of objects, which we recombined into our scene-graph based captions with complexity levels ranging from 3 to 9. 1282 1283 Finally, we used DaLL-E 3 to generate corresponding im-1284 ages for these newly composed captions.

1285The fine-tuning configurations were consistent with1286those used in the self-improving setup (Appendix E.1.3). To1287accommodate the reduced dataset size, the maximum train-1288ing steps were set to 1000.

1289As a baseline, we randomly selected 778 images from129010k GENERATE ANY SCENE-generated images, using cap-1291tions produced by GENERATE ANY SCENE. This ensured1292a controlled comparison between the targeted and random1293fine-tuning strategies.

G. Details of Application 3: Generated content detector (Section 3) 1295

1296

G.1. Experiment details

In this section, our goal is to validate that the more diverse
captions generated by GENERATE ANY SCENE can com-
plement existing datasets, which are predominantly com-
posed of real-world images paired with captions. By do-
ing so, we aim to train AI-generated content detectors to
achieve greater robustness.1297
12981302

Dataset preparation We conducted comparative exper-1303 iments between captions generated by GENERATE ANY 1304 SCENE and entries from the D^3 dataset. From the D^3 1305 dataset, we randomly sampled 10k entries, each including 1306 a caption, a link to a real image, and an image generated by 1307 SD v1.4. Due to some broken links, we successfully down-1308 loaded 8.5k real images and retained 10k SD v1.4-generated 1309 images. We also used SD v1.4 to generate images based on 1310 10k GENERATE ANY SCENE captions. 1311

We varied the training data sizes based on the sampled 1312 dataset. Specifically, we sampled N real images from the 1313 10k D^3 real images. For synthetic data, we compared N 1314 samples exclusively from D^3 with a mixed set of N/2 sam-1315 ples from 10k GENERATE ANY SCENE images and N/2 1316 sampled from D^3 , ensuring a total of N synthetic samples. 1317 Combined, this resulted in 2N training images. We tested 1318 2N across various sizes, ranging from 2k to 10k. 1319

Detector architecture and training We employed ViT-1320 T [62] and ResNet-18 [18] as backbones for the detection 1321 models. Their pretrained parameters on ImageNet-21k were 1322 frozen, and the final classification head was replaced with 1323 a linear layer using a sigmoid activation function to pre-1324 dict the probability of an image being AI-generated. Dur-1325 ing training, We used Binary Cross-Entropy (BCE) as the 1326 loss function, and the AdamW optimizer was applied with a 1327 learning rate of $2e^{-3}$. Training was conducted with a batch 1328 size of 256 for up to 50 epochs, with early stopping trig-1329 gered after six epochs of no improvement in validation per-1330 formance. 1331

TestingTo evaluate the performance of models trained1332with varying dataset sizes and synthetic data combina-
tions, we tested them on both GenImage and GENERATE1333ANY SCENE datasets to assess their in-domain and out-of-
domain performance under different settings.1332

For GenImage, we used validation data from four mod-
els: SD v1.4, SD v1.5, MidJourney, and VQDM. Each val-
idation set contained 8k real images and 8k generated im-
ages. For GENERATE ANY SCENE, we sampled 10k real
images from CC3M and paired them with 10k generated1337
1338

1342	images from each of the following models: Stable Diffu-
1343	sion v2-1, PixArt- α , Stable Diffusion 3 Medium, and Play-
1344	ground v2.5. This created distinct test sets for evaluating
1345	model performance across different synthetic data sources.

1346 G.2. Results

Table 10 and Table 9 evaluate the performance of ResNet-1347 18 and ViT-T detection backbones trained on datasets of 1348 varying sizes and compositions across in-domain (same 1349 model and cross-model) and out-of-domain settings. While 1350 models trained with D^3 and GENERATE ANY SCENE oc-1351 casionally underperform compared to those trained solely 1352 on D^3 in the in-domain same-model setting, they exhibit 1353 significant advantages in both in-domain cross-model and 1354 out-of-domain evaluations. These results demonstrate that 1355 incorporating our data (GENERATE ANY SCENE) into the 1356 1357 training process enhances the detector's robustness. By supplementing existing datasets with GENERATE ANY SCENE 1358 1359 under the same training configurations and dataset sizes, detectors achieve stronger cross-model and cross-dataset ca-1360 pabilities, highlighting improved generalizability to diverse 1361 generative models and datasets. 1362

Detector	Data Scale	SDv1.4 (In-domain, same model)		SDv2.1		Pixart	Pixart- α		SDv3-medium		Playground v2.5		Average (In-domain, cross model)	
	(2N)	D^3 + Ours	D^3	D^3 + Ours	D^3	D^3 + Ours	D^3	D^3 + Ours	D^3	D^3 + Ours	D^3	D^3 + Ours	D^3	
	2k	0.6561	0.6663	0.7682	0.6750	0.7379	0.606	0.7509	0.6724	0.7380	0.5939	0.7488	0.6368	
	4k	0.6751	0.6812	0.7624	0.6853	0.7328	0.6494	0.7576	0.7028	0.7208	0.6163	0.7434	0.6635	
Resnet-18	6k	0.6780	0.6995	0.7886	0.6870	0.7493	0.6586	0.7768	0.7285	0.7349	0.6335	0.7624	0.6769	
	8k	0.6828	0.6964	0.7710	0.6741	0.7454	0.6418	0.7785	0.7186	0.7215	0.6033	0.7541	0.6595	
	10k	0.6830	0.6957	0.7807	0.6897	0.7483	0.6682	0.7781	0.7326	0.7300	0.6229	0.7593	0.6784	
	2k	0.6759	0.6672	0.7550	0.6827	0.7585	0.6758	0.7473	0.6941	0.7327	0.6106	0.7484	0.6658	
	4k	0.6878	0.6871	0.7576	0.7000	0.7605	0.7071	0.7549	0.7217	0.7221	0.6144	0.7488	0.6858	
ViT-T	6k	0.6898	0.6891	0.7663	0.6962	0.7666	0.7164	0.7629	0.7238	0.7303	0.6134	0.7565	0.6875	
	8k	0.6962	0.6974	0.7655	0.6894	0.7712	0.7253	0.7653	0.7253	0.7381	0.6344	0.7600	0.6936	
	10k	0.6986	0.6984	0.7828	0.6960	0.7777	0.7275	0.7786	0.7334	0.7330	0.6293	0.7680	0.6966	

Table 9. F1-Score Comparison of ResNet-18 and ViT-T Detectors Trained with D^3 and D^3 + GENERATE ANY SCENE Across In-Domain Settings

Detector	Data Scale (2N)	SDv1.5		VQDM		Midjourney		Average (Out-of-domain)	
		D^3 + Ours	D^3	D^3 + Ours	D^3	D^3 + Ours	D^3	D^3 + Ours	D^3
Resnet-18	2k	0.6515	0.6591	0.5629	0.5285	0.5803	0.5647	0.5982	0.5841
	4k	0.6709	0.6817	0.5693	0.5428	0.6016	0.5941	0.6139	0.6062
	6k	0.6750	0.6963	0.5724	0.5327	0.6084	0.6072	0.6186	0.6121
	8k	0.6792	0.6965	0.5716	0.5282	0.6097	0.5873	0.6202	0.6040
	10k	0.6814	0.6955	0.5812	0.5454	0.6109	0.6040	0.6245	0.6150
ViT-T	2k	0.6755	0.6685	0.5443	0.4966	0.6207	0.6066	0.6135	0.5906
	4k	0.6845	0.6865	0.5591	0.4971	0.6416	0.6149	0.6284	0.5995
	6k	0.6900	0.6890	0.5580	0.4948	0.6455	0.6259	0.6313	0.6032
	8k	0.6940	0.6969	0.5553	0.4962	0.6495	0.6387	0.6329	0.6106
	10k	0.6961	0.6988	0.5499	0.4975	0.6447	0.6358	0.6302	0.6107

Table 10. F1-Score Comparison of ResNet-18 and ViT-T Detectors Trained with D^3 and D^3 + GENERATE ANY SCENE Across Out-of-Domain Settings