
Copy2Prompt: Copyright-Compliant Image Synthesis via Text-Mediated Generation

Lucas Lee

Department of Data Science
UCSD
San Diego, CA, 92092
l2lee@ucsd.edu

Anandhini Rajendran

Department of Computer Science
UCSD
San Diego, CA, 92092
anrajendran@ucsd.edu

David Lurie

Department of Data Science
UCSD
San Diego, CA, 92092
dlurie@ucsd.edu

Abstract

1 We introduce Copy2Prompt, a novel framework for copyright-compliant image-to-
2 image synthesis that resolves the fundamental tension between semantic fidelity
3 and ethical content generation by operating entirely in the textual domain. Our
4 three-stage pipeline extracts detailed descriptions from reference images using
5 vision-language models like LLaVA, refines these descriptions through large lan-
6 guage models or RL agents, and synthesizes variants using diffusion models without
7 direct pixel reuse. Evaluation across multiple backends (DALL-E 3, Stable Diffu-
8 sion XL) demonstrates it is possible to utilize just textual description to generate
9 faithful image while eliminating copyright concerns. Copy2Prompt represents a
10 principled solution for ethical generative AI that balances creative flexibility with
11 legal compliance through text-mediated synthesis.

1 Introduction

13 Diffusion models have transformed both text-to-image and image-to-image (I2I) synthesis, achieving
14 photorealism and stylistic flexibility once reserved for GANs or classical renderers [1, 2]. In the
15 I2I setting, a user supplies a reference picture and requests an edited or stylistically altered version.
16 State-of-the-art pipelines typically fall into two families:

- 17 1. **Pixel/latent inversion.** Methods such as Stable Diffusion XL (SDXL) [3] invert the reference
18 image into the diffusion latent space, then sample variants conditioned on that latent. While
19 they preserve fidelity, they often reproduce protected pixels verbatim, raising copyright and
20 ethical concerns [4].
- 21 2. **Conditional control.** Approaches like ControlNet inject edge, pose or depth hints into the
22 denoising process [5]. They guarantee a geometric match but cannot retain fine appearance
23 details without, again, copying large latent fragments of the original.

24 Thus, increasing noise during inversion avoids plagiarism but drifts semantically; decreasing noise
25 retains detail but risks near-duplicates. The community therefore lacks a principled mechanism to
26 balance *semantic fidelity* and *copyright safety*.

27 **Text as a safe, controllable bridge.** Recent vision–language models—e.g. BLIP-2 [6] and LLaVA
28 [7]—can translate an image into a rich, grounded description. Concurrently, text-to-image back-ends
29 such as SDXL [3] and single-step distillations like ADD/SD-Turbo [8] accept hundreds of tokens
30 and expose fine-grained prompt weighting [9, 10]. These developments suggest a new paradigm:
31 *convert the reference image to text, edit or weight that text, then re-synthesise from scratch*. Because
32 generation is conditioned solely on words, no protected pixels are reused.

33 Thus, we introduce Copy2Prompt, a text-mediated I2I pipeline simultaneously offers fine image
34 semantic control and strong copyright compliance. The pipeline can be visualized in Figure 1.

35 2 Background

36 Our Copy2Prompt pipeline builds on recent advances in multimodal learning that enable high-fidelity
37 visual-to-text translation through CLIP-based encoders and instruction-tuned language models. On
38 the generation side, we leverage state-of-the-art diffusion and transformer-based models—particularly
39 DALL-E 3 and SDXL—augmented with Compel’s prompt weighting mechanisms to regain semantic
40 and stylistic alignment with the original image. Finally, we enhance prompts using LLMs, and also
41 experiment with reinforcement learning using CLIP-based scoring to evaluate prompt-generated
42 images. Together, these components allow for controlled, copyright-compliant image regeneration
43 guided by precise textual representations.

44 Visual-to-Textual Translation

45 The landscape of image-to-text translation has evolved significantly with the emergence of various
46 multimodal architectures. Contrastive Language-Image Pretraining (CLIP) learns image represen-
47 tations that capture semantics and style, allowing it to serve as the basis for prompting image
48 regeneration that is similar to the original image [11]. Recent work, however, has demonstrated
49 that CLIP models inherently possess the ability to convert images into text prompts through linear
50 projection matrices calculated in closed form [12]. Liu et al.’s visual instruction tuning established
51 that language models can process visual information through frozen CLIP encoders, achieving 85.1%
52 relative performance on reasoning tasks [13]. This foundation enables pixel-accurate visual grounding
53 in OMG-LLaVA (Zhang et al.) [14], which integrates universal segmentation with large language
54 models (LLMs).

55 Text-to-Image Generation

56 Contemporary work in text-to-image generation has evolved from more rigid models to complex multi-
57 component pipelines that combine specialized generators, adapters, and upscaling modules [14] [15].
58 Diffusion models remain dominant in this space, with some providing highly controllable frame-
59 works [16] [17].

60 Methods to optimize these frameworks, like policy optimization, reduce computation and improve
61 quality by training ensemble reward models that predict image quality directly from prompt-workflow
62 combinations [18]. Prompt engineering has also been explored, revealing that different timesteps in
63 diffusion processes exhibit varying sensitivity to conditioning, with noisy later steps showing greater
64 responsiveness to prompt modifications [19].

65 **Compel Integration for Prompt Weighting** Compel [20] is a library that enables direct manipulation
66 of text embedding tensors in diffusion models. Compel parses weighted prompt syntax and modifies
67 the corresponding CLIP text encoder embeddings before they condition the diffusion process. The
68 library applies scalar multipliers to specific token embeddings: + operators scale embeddings by 1.1
69 \times while - operators apply $0.9 \times$ scaling, with compound operations enabling fine-grained attention
70 control (e.g., $++$ results in $1.1^2 = 1.21$). The .and() function concatenates separately weighted
71 embedding segments, allowing heterogeneous conditioning where different prompt components
72 receive distinct emphasis levels during the denoising process.

73 **DALL-E 3** DALL-E 3’s native ChatGPT integration automatically refines prompts, expanding
74 descriptions into detailed generation instructions that complement LLaVA’s annotations. The model’s
75 enhanced nuance understanding enables more accurate translation of textual descriptions into visual
76 outputs. Additionally, DALL-E 3’s sophisticated content filtering and copyright compliance measures
77 align with Copy2Prompt’s objectives, providing built-in safeguards against copyrighted material
78 generation while maintaining creative flexibility.

79 3 Methodology

80 We propose Copy2Prompt, a framework designed to generate high-fidelity image variations by
81 converting source images into detailed textual prompts, which are then used to guide advanced text-
82 to-image diffusion models. The pipeline integrates vision-language models and prompt engineering

83 to achieve copyright-compliant, semantically rich image synthesis, with the latter vision-language
84 diffusion model never having no pixel-level knowledge of the source image, thereby avoiding
85 copyright claims.

86 **3.1 Pipeline Overview**

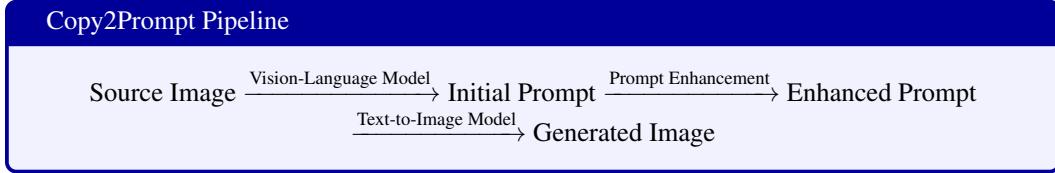


Figure 1: Overview of the Copy2Prompt pipeline for generating image variations through text-mediated synthesis.

87 As seen in Figure 1, Copy2Prompt works in three steps.

- 88 1. Image understanding (vision language model):
 - 89 • Input: A source image
 - 90 • Output: A detailed textual prompt describing the visual content, objects, and relationships present in the image.
- 92 2. Prompt Enhancement:
 - 93 • Input: Output from previous step
 - 94 • Output: Enhanced textual prompt with improved coherence, tailoring the description
 - 95 for downstream image synthesis
- 96 3. Image Synthesis (text-to-image diffusion model):
 - 97 • Input: Output from previous step
 - 98 • Output: Generated image, replicating the source image through pure textual description

99 **3.2 Image Understanding**

100 The image understanding component transforms visual information into rich textual representations
101 that capture both semantic content and spatial relationships. We employ vision-language models,
102 specifically LLaVA, to generate comprehensive descriptions that go beyond simple object detection,
103 incorporating contextual details, visual attributes, and compositional elements essential for high-
104 fidelity image reconstruction through text-to-image synthesis.

105 LLaVA yields multi-sentence annotations that exceed traditional captioning approaches in both
106 depth and granularity. The architecture converts images into text through a two-stage process: first,
107 CLIP's vision encoder extracts visual features from the input image, then these features are passed
108 to a modified LLaMA language model, which generates comprehensive descriptive captions. The
109 complete LLaVA processing pipeline is illustrated in Figure 2.

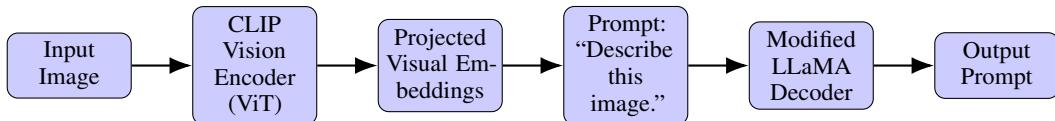


Figure 2: LLaVA architecture for converting input images to detailed textual descriptions through CLIP visual encoding and LLaMA language generation.

110 **3.3 Prompt Enhancement**

111 The initial prompt generated by LLaVA often requires strategic refinement to optimize compatibility
112 with text-to-image generation models. We employ large language models (ChatGPT) to enhance

113 these descriptions. The system prompt used by ChatGPT to format the detailed prompt is shown in
114 Appendix A.

115 The prompt enhancement process addresses several key limitations of raw vision-language outputs:
116 improving linguistic coherence, expanding contextual details, and restructuring descriptions to align
117 with the inductive biases of diffusion models. This refinement ensures that critical visual elements
118 are preserved while adapting the textual representation for optimal downstream synthesis.

119 **3.4 Image Generation**

120 The final stage of our pipeline leverages state-of-the-art text-to-image diffusion models to synthesize
121 high-fidelity images from the enhanced textual descriptions. We primarily employ Stable Diffusion XL
122 (SDXL) as our generation backbone, though our framework supports multiple diffusion architectures
123 to accommodate different quality-speed trade-offs.

124 SDXL and diffusion models in general can only take in a textual inputs of 77 tokens. This is because
125 CLIP, the underlying text encoder, is limited to processing 77 tokens per prompt.

126 To address this limitation, we implement two complementary strategies to maximize prompt utilization
127 and preserve the richness of information generated during the enhancement stage.

128 First, we integrate Compel with SDXL to enable advanced prompt processing capabilities, includ-
129 ing weighted attention mechanisms, prompt segmentation for fine-grained control, and intelligent
130 chunking with padding to retain maximum textual information within model constraints.

131 Second, we support DALL-E 3 as an alternative generator, which allows for prompts up to 4000
132 tokens — significantly expanding the context window compared to base SDXL for the ability to read
133 in detailed prompts.

134 **3.5 Reinforcement as an Alternative for Prompt Enhancement**

135 Instead of relying on fixed or manually engineered prompts and having an LLM revise prompts, we
136 formulate the problem as a prompt optimization for image generation as an RL problem, where the
137 agent iteratively modifies text prompts to maximize visual similarity between the target image and
138 generated image. The agent interacts with the text-to-image generation and is guided by a CLIP-based
139 reward. This approach makes no assumptions on the image generator, making it scalable for larger
140 pipelines.

141 **RL Formulation:** The agent interacts with the Stable Diffusion XL generator by sampling prompt
142 mutations (actions) to generate images from text (states), aiming to maximize a CLIP-based simi-
143 larity score (reward) with a reference image. The policy randomly selects stylistic, technical, and
144 compositional modifiers to evolve prompts toward higher visual-semantic alignment.

145 **Optimization Loop** The agent iteratively changes the best prompt, generates an image via the SDXL
146 pipeline, evaluates it with CLIP similarity to the reference, and updates the best prompt and image if
147 the reward improves.

148 **4 Evaluation**

149 In this section, we outline both the similarity metrics used to assess image quality and the experimental
150 results that demonstrate the effectiveness of our pipeline. This comprehensive evaluation provides
151 insight into how well our approach preserves visual fidelity and semantic alignment, as well as
152 parameters and different components affect the generated images.

153 **4.1 Similarity Metrics:**

154 To assess the quality of the generated images, we evaluate their similarity to the original images using
155 a combination of perceptual and semantic similarity metrics:

- 156 • **Structural Similarity Index (SSIM)** measures the perceptual similarity in structure and
157 texture between two images, reflecting how well the layout and content are preserved.

- 158 • **CLIP Score** evaluates semantic similarity between images and text, capturing high-level
 159 content alignment and conceptual consistency.
 160 • **Learned Perceptual Image Patch Similarity (LPIPS)** quantifies perceptual distance using
 161 deep features, providing a measure of human-perceived visual similarity.

162 **4.2 Experiments and Results**

163 **4.2.1 RL Agent for Prompt Enhancement**

164 We use different type of prompt modifiers in the RL agent to evaluate the impact on the generated
 165 image. The prompt modifiers were categorized from three groups:

- 166 • **Quality Enhancers:** These include descriptors that improve image fidelity and realism,
 167 such as “hyperrealistic” and “8k resolution.”
 168 • **Technical Specs:** Modifiers reflecting photographic parameters, for example, “shot on
 169 Canon EOS R5” and “f/1.4 aperture.”
 170 • **Composition Rules:** Terms guiding the visual arrangement, like “rule of thirds” and
 171 “symmetrical composition.”

172 The RL agent employs an iterative optimization strategy, sampling from predefined modifier pools
 173 to progressively enhance prompt effectiveness. We demonstrate this process using a representative
 174 example with the base prompt: “*Sunset at the beach with a group of people.*” as shown in Table 1
 175 and Figure 3

Table 1: RL Agent Optimization Trajectory

Iteration	Enhanced Prompt	CLIP Score	Key Modifiers
Base	“Sunset at the beach with a group of people”	—	—
1	Base + “intricate details, shot on Canon EOS R5, balanced composition”	0.7515	Technical + Composition
2	Base + “ultra realistic, long exposure, golden ratio”	0.7129	Quality + Technical + Composition
3	Base + “fine art quality, film grain texture, leading lines”	0.7593	Quality + Technical + Composition
4	Base + “award winning, bokeh background, diagonal composition”	0.8228	Quality + Technical + Composition
5	Base + “professional photography, f/1.4 aperture, negative space”	0.7905	Quality + Technical + Composition

176 **Convergence Pattern:** The optimization process demonstrates non-monotonic convergence, with
 177 peak performance achieved at iteration 4 (CLIP score: 0.8228), followed by a slight decline in
 178 iteration 5.

179 **Optimal Modifier Combination:** The highest-performing configuration combined quality descriptors (“award winning”), technical specifications (“bokeh background”), and compositional guidelines (“diagonal composition”).

182 **Implications:** This case study demonstrates that RL agents can systematically transform simple
 183 base prompts into sophisticated generation instructions through strategic modifier addition. The 9.5%
 184 improvement from baseline to optimal configuration ($0.7515 \rightarrow 0.8228$) validates the effectiveness
 185 of compositional and stylistic enhancement for image-text alignment tasks. The non-monotonic
 186 optimization trajectory highlights the importance of allowing sufficient exploration in the RL training
 187 process, as premature convergence could miss optimal modifier combinations.

188 Qualitatively, the source and generated image pair in Figure 1 are semantically similar to humans.

189 Figures 4 and Figure 5 show additional experiments for the RL agent. Qualitatively, the lake images
 190 are compositionally similar, where as the cat images are not. This is in contrast to what CLIP score



Figure 3: Transformation from input image (left) to output image (center) generated by the RL agent, with the final CLIP similarity score and prompt shown on the right.

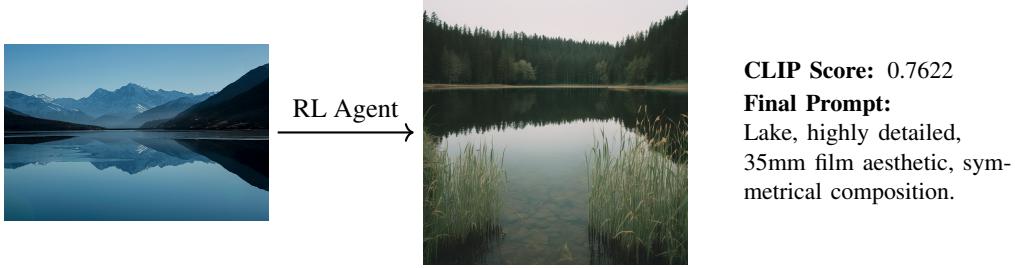


Figure 4: Transformation from input image (left) to output image (center) generated by the RL agent, with the final CLIP similarity score and prompt shown on the right.

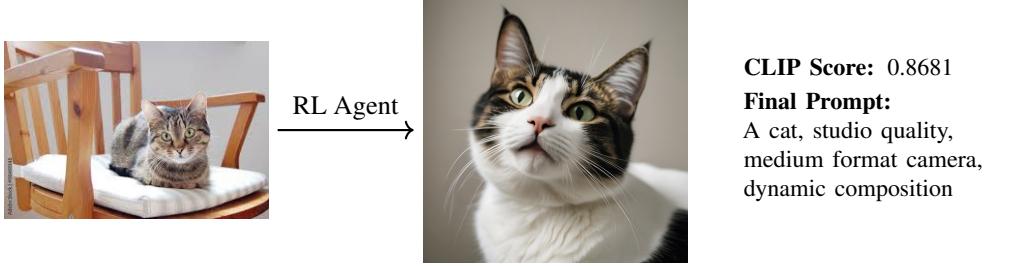


Figure 5: Transformation from input image (left) to output image (center) generated by the RL agent, with the final CLIP similarity score and prompt shown on the right.

191 is suggesting. This suggests that CLIP score may value color accuracy more than compositional
192 accuracy, and this prioritization does not always align with human judgement.

193 4.2.2 GPT Model and Prompt Length into SDXL + Compel

194 We evaluate the effectiveness of Copy2Prompt on a subset of the Unsplash dataset ($N = 15$) containing
195 photographs akin to background wallpapers. Our experimental setup in Table 2 systematically varies
196 the ChatGPT model version (GPT-4.0 vs GPT-4.1) and target prompt length (77, 200, 500, and 1000
197 words) to assess their impact on generation quality.

198 For each image, we generate initial descriptions using LLaVA, enhance these prompts using ChatGPT
199 with specific length constraints, then feed the enhanced prompts to SDXL for image generation. We
200 evaluate results using SSIM (structural similarity), LPIPS (perceptual distance), CLIP score (semantic
201 alignment), and actual output prompt length.

202 Initially, we used GPT-4.0 with explicit word count constraints but observed significant length
203 mismatches in actual output, particularly for longer targets (500-1000 words), prompting a switch to
204 GPT-4.1 which demonstrated better constraint adherence. The results reveal that longer enhanced
205 prompts do not consistently improve generation quality: GPT-4.0 shows decreasing SSIM scores

Table 2: Comparison of model performance on varying Prompt Enhancement Model and Prompt Length

Model/Technique	Mean SSIM \uparrow	Mean LPIPS \downarrow	Mean CLIP Score \uparrow	Mean output length from ChatGPT
Copy2Prompt 4o (77)	0.19	0.67	0.80	72
Copy2Prompt 4o (200)	0.19	0.68	0.79	180
Copy2Prompt 4o (500)	0.17	0.66	0.79	383
Copy2Prompt 4o (1000)	0.16	0.68	0.79	211
Copy2Prompt 4.1 (77)	0.18	0.68	0.80	73
Copy2Prompt 4.1 (200)	0.20	0.67	0.81	182
Copy2Prompt 4.1 (500)	0.18	0.67	0.81	450
Copy2Prompt 4.1 (1000)	0.18	0.70	0.80	890

from 0.19 (77 words) to 0.16 (1000 words) with stable CLIP scores (0.79-0.80), while GPT-4.1 maintains consistent performance across lengths (SSIM: 0.18-0.20, CLIP: 0.80-0.81). Notably, both models achieve similar performance with concise 77-word prompts, suggesting well-crafted shorter descriptions can be as effective as verbose ones.

4.2.3 LLaVa Prompt Length + ChatGPT Prompt Length:

Table 3: Comparison of model performance on varying LlaVA Output Length

Model/Technique	Mean SSIM \uparrow	Mean LPIPS \downarrow	Mean CLIP Score \uparrow	Mean output length from ChatGPT
Copy2Prompt 4.1 (77)	0.18	0.68	0.79	75
Copy2Prompt 4.1 (200)	0.16	0.69	0.81	186
Copy2Prompt 4.1 (500)	0.17	0.67	0.81	460
Copy2Prompt 4.1 (1000)	0.17	0.70	0.83	820

To investigate whether richer initial descriptions improve generation quality, we configured LLaVA with extended token limits (1800-2400 tokens) compared to (600-1000) from Table 2 and repeated the experiment using the same dataset ($N = 15$) and GPT-4.1 enhancement with identical target lengths. The results in Table 3 compared to those in Table 2 show modest improvements in semantic alignment, where CLIP scores demonstrate a clearer upward trend from 0.79 (77 words) to 0.83 (1000 words). On the other hand, SSIM and LPIPS values remain largely unchanged. While longer LLaVA outputs slightly improve semantic metrics, the effect size is minimal and unlikely to significantly impact practical image generation quality.

4.2.4 Image Generators + ChatGPT Prompt Length

Table 4: Comparison of model performance across different image generators and prompt length

Model/Technique	Mean SSIM \uparrow	Mean LPIPS \downarrow	Mean CLIP Score \uparrow
SDXL Source Image to Image Baseline	0.62	0.16	0.92
Source image → ChatGPT → Dalle-3 Baseline	0.20	0.67	0.82
BLIP-2 → SDXL Baseline	0.11	0.75	0.66
Copy2Prompt (Dalle-3, 77)	0.20	0.72	0.81
Copy2Prompt (Dalle-3, 200)	0.19	0.70	0.78
Copy2Prompt (Dalle-3, 500)	0.22	0.69	0.79
Copy2Prompt (Dalle-3, 1000)	0.21	0.70	0.77
Copy2Prompt (SDXL+C, 77)	0.20	0.68	0.79
Copy2Prompt (SDXL+C, 200)	0.17	0.69	0.77
Copy2Prompt (SDXL+C, 500)	0.16	0.67	0.78
Copy2Prompt (SDXL+C, 1000)	0.16	0.66	0.79

220 To evaluate different image generation backends, we compared Copy2Prompt using DALL-E 3 and
221 SDXL (with Compel) against three baselines: SDXL image-to-image (direct pixel-level generation),
222 ChatGPT→DALL-E 3 (image-to-text-to-image), and BLIP-2→SDXL (captioning baseline).

223 As shown in Table 4, SDXL image-to-image, as expected, achieves superior performance (CLIP:
224 0.92, SSIM: 0.62), reflecting the advantages of direct pixel manipulation. Among text-mediated
225 approaches, Copy2Prompt variants consistently outperform the BLIP-2 captioning baseline and
226 perform comparably to the ChatGPT baseline. The similarity between ChatGPT→DALL-E 3 and
227 Copy2Prompt performance suggests that LLaVA (llava-hf/llava-v1.6-mistral-7b-hf), despite being de-
228 signed for multi-granular image description, may lack sufficient detail for optimal prompt generation.
229 Notably, the ChatGPT baseline can be viewed as an alternative Copy2Prompt implementation, as
230 both perform image-to-text-to-image translation.

231 Our results reveals distinct backend-specific responses to prompt length variations. SDXL+Compel
232 shows degrading structural similarity as prompt length increases (SSIM: 0.22→0.16) while main-
233 taining stable perceptual metrics (LPIPS, CLIP). Conversely, DALL-E 3 preserves structural and
234 perceptual consistency across lengths but exhibits declining semantic alignment (CLIP: 0.81→0.77).
235 These contrasting behaviors demonstrate that optimal prompt strategies are model-dependent, with
236 prompt quality and relevance proving more critical than length alone.

237 A full run of Copy2Prompt can be seen in Appendix B.

238 5 Conclusion

239 5.1 Discussion

240 Copy2Prompt demonstrates an image-text-image pipeline that enhances image generation through
241 prompt optimization and model selection. The RL agent successfully alters prompts by adding details
242 and wording that results in generated images that have greater semantic similarity to the original
243 image. Furthermore, our experiments show that increasing the length of the image-understanding
244 prompt improves image generation quality. However, increasing the prompt enhancement prompt
245 length does not always result in improved quality, and the effect of the prompt length depends on the
246 image generation model used. When comparing candidate image generators, DALLE-3 performs
247 marginally better than SDXL with Compel in our Copy2Prompt framework. Our results also highlight
248 the limitations of captioning models like LLaVA for comprehensive image descriptions, where their
249 textual descriptions appear no better than general purpose multimodal models like ChatGPT 4o.

250 Overall, we believe that it's possible to generate semantically similar photos to the human eyes
251 without using any pixel level information based on the experiments. However, significant effort has to
252 go into the prompting of the image understanding model, the prompt enhancer model, or the compute
253 spent on training RL agents.

254 5.2 Future Work

255 Future research directions include:

- 256 1. **Enhanced RL Integration:** Extending RL optimization to the vision-language understand-
257 ing stage and integrating with LLMs for dynamic prompt generation rather than sampling
258 from curated modifier lists
- 259 2. **Advanced Description Strategies:** Since longer descriptions did not improve fidelity,
260 exploring ensemble methods where LLaVA generates multiple complementary descriptions,
261 hierarchical frameworks separating structural and semantic elements, and more sophisticated
262 vision-language models, like OMG-LlaVA
- 263 3. **Copyright Validation:** Systematic evaluation on real copyright infringement datasets to
264 establish empirical guidelines for transformative use and validate legal compliance.

265 References

- 266 [1] J. Ho, A. Jain, and P. Abbeel, “Denoising diffusion probabilistic models,” *Advances in Neural*
267 *Information Processing Systems*, vol. 33, pp. 6840–6851, 2020.

- 268 [2] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, “High-resolution image
269 synthesis with latent diffusion models,” in *CVPR*, 2022, pp. 10 684–10 695.
- 270 [3] A. Podell and S. AI, “Stable diffusion xl,” <https://stability.ai/blog/stable-diffusion-xl>, 2023.
- 271 [4] N. Yu, C. Meng, M. Fritz *et al.*, “Synthetic imagery can be traced back to its generative source,”
272 *Nature*, vol. 620, pp. 104–111, 2023.
- 273 [5] L. Zhang and M. Agrawala, “Adding conditional control to text-to-image diffusion models,”
274 *arXiv preprint arXiv:2302.05543*, 2023.
- 275 [6] J. Li, D. Li, C. Xiong, and S. Hoi, “Blip-2: Bootstrapping language-image pre-training with
276 frozen image encoders and large language models,” *ICLR*, 2023.
- 277 [7] H. Liu, Q. Ji *et al.*, “Visual instruction tuning,” *arXiv preprint arXiv:2304.08485*, 2023.
- 278 [8] A. Sauer, N. Frühstück, M. Sun, and A. Geiger, “Adversarial diffusion distillation,” *ICLR*, 2023.
- 279 [9] A. Hertz, Y. Bahat, R. Giryes, A. A. Efros, and E. Shechtman, “Prompt-to-prompt image editing
280 with cross attention control,” in *SIGGRAPH Asia*, 2022.
- 281 [10] S. Mahajan, T. Rahman, K. M. Yi, and L. Sigal, “Prompting hard or hardly
282 prompting: Prompt inversion for text-to-image diffusion models,” 2023. [Online]. Available:
283 <https://arxiv.org/abs/2312.12416>
- 284 [11] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell,
285 P. Mishkin, J. Clark, G. Krueger, and I. Sutskever, “Learning transferable visual models from
286 natural language supervision,” 2021. [Online]. Available: <https://arxiv.org/abs/2103.00020>
- 287 [12] Y. Ding, C. Tian, H. Ding, and L. Liu, “The clip model is secretly an image-to-prompt
288 converter,” 2024. [Online]. Available: <https://arxiv.org/abs/2305.12716>
- 289 [13] H. Liu, C. Li, Q. Wu, and Y. J. Lee, “Visual instruction tuning,” 2023. [Online]. Available:
290 <https://arxiv.org/abs/2304.08485>
- 291 [14] C. Zhang, C. Zhang, M. Zhang, I. S. Kweon, and J. Kim, “Text-to-image diffusion models in
292 generative ai: A survey,” 2024. [Online]. Available: <https://arxiv.org/abs/2303.07909>
- 293 [15] C. Saharia, W. Chan, S. Saxena, L. Li, J. Whang, E. Denton, S. K. S. Ghasemipour, B. K. Ayan,
294 S. S. Mahdavi, R. G. Lopes, T. Salimans, J. Ho, D. J. Fleet, and M. Norouzi, “Photorealistic
295 text-to-image diffusion models with deep language understanding,” 2022. [Online]. Available:
296 <https://arxiv.org/abs/2205.11487>
- 297 [16] S. Hartwig, D. Engel, L. Sick, H. Kniesel, T. Payer, P. Poonam, M. Glöckler, A. Bäuerle,
298 and T. Ropinski, “A survey on quality metrics for text-to-image generation,” 2025. [Online].
299 Available: <https://arxiv.org/abs/2403.11821>
- 300 [17] X. Zhang, H. Zhou, H. Qin, X. Lu, J. Yan, G. Wang, Z. Chen, and Y. Liu, “Enabling versatile
301 controls for video diffusion models,” 2025. [Online]. Available: <https://arxiv.org/abs/2503.16983>
- 302 [18] U. Gadot, R. Gal, Y. Ziser, G. Chechik, and S. Mannor, “Policy optimized text-to-image
303 pipeline design,” 2025. [Online]. Available: <https://arxiv.org/abs/2505.21478>
- 304 [19] S. Mahajan, T. Rahman, K. M. Yi, and L. Sigal, “Prompting hard or hardly
305 prompting: Prompt inversion for text-to-image diffusion models,” 2023. [Online]. Available:
306 <https://arxiv.org/abs/2312.12416>
- 307 [20] damian0815, “Compel,” 2025. [Online]. Available: <https://github.com/damian0815/compel>

308 **Appendix**

309 **A System Prompts**

310 This section details the system prompts used for different components of the Copy2Prompt pipeline.

311 **A.1 LLaVA Enhancement Prompt**

312 The following system prompt was used to enhance initial LLaVA-generated descriptions:

LLaVA Enhancement System Prompt

You are an expert at converting detailed image descriptions into concise, effective prompts for Stable Diffusion XL.

Your task is to transform a long, detailed description into a prompt of 20% more or less than EXACTLY {target_length} words.

Key requirements:

1. Preserve the most visually important elements
2. Use comma-separated phrases for clarity
3. Front-load the most important elements (subject, style, composition)
4. Remove redundant or verbose descriptions
5. Keep specific visual details that define the image
6. MUST be {target_length} words or less

313

314 **A.2 ChatGPT Base Prompt**

315 The base system prompt for ChatGPT as the prompt enhancer:

ChatGPT Base System Prompt

You are an expert at converting detailed image descriptions into concise, effective prompts for Stable Diffusion XL.

Your task is to transform a long, detailed description into a prompt of 20% more or less than EXACTLY {target_length} words.

Key requirements:

1. Preserve the most visually important elements
2. Use comma-separated phrases for clarity
3. Front-load the most important elements (subject, style, composition)
4. Remove redundant or verbose descriptions
5. Keep specific visual details that define the image
6. MUST be {target_length} words or less

316

317 **A.3 ChatGPT Style Instructions**

318 Additional style instructions provided to ChatGPT for structured prompt generation:

ChatGPT Style Instructions

Include: main subject with key features, composition, artistic style, atmosphere, important details

Structure: [detailed subject], [composition], [style/medium], [atmosphere/mood], [key details]

319

320 **A.4 ChatGPT Fallback Prompt**

321 When the ChatGPT formatter produces output that deviates from the target length by more than 20%,
322 the following fallback prompt is used:

ChatGPT Length Correction Fallback

Convert the description to EXACTLY {target_length} words.
Count every word. Be extremely concise.

323

324 **B Case Study: Complete Copy2Prompt Pipeline**

325 This section demonstrates a complete Copy2Prompt execution using GPT-4.1 for prompt formatting
326 with a target length of 77 words. The LLaVA model was configured to generate descriptions between
327 800-1200 tokens to ensure comprehensive image analysis, and SDXL with Compel is used to generate
328 the final image.

329 **B.1 Pipeline Execution**

330 **Stage 1: Source Image Analysis**

331 Figure 6 shows the input image from the Unsplash dataset, depicting a winter forest scene with a
332 cross-country skier.

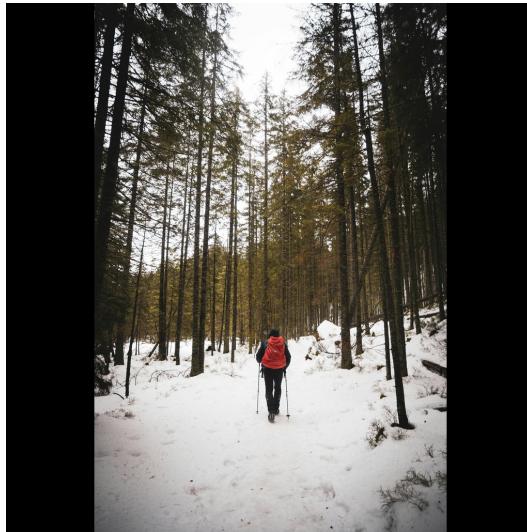


Figure 6: Source image: Winter forest scene with cross-country skier

333 **Stage 2: LLaVA Description Generation**

334 The vision-language model produced a detailed description scene composition, lighting conditions,
335 and photographic elements. Key extracted features include:

- 336 • *Subject*: Lone figure in red jacket with ski poles on snow-covered path
- 337 • *Environment*: Dense forest with snow-dusted coniferous and deciduous trees
- 338 • *Composition*: Off-center positioning following rule of thirds, blurred background for depth
- 339 • *Lighting & Color*: Soft diffused lighting, muted monochromatic palette with red accent
- 340 • *Style*: Professional landscape photography with attention to texture and detail

341 The full prompt below:

342 “The image captures a serene winter scene in a forest. The main subject is a lone
343 figure, clad in a vibrant red jacket and black pants, walking away from the viewer
344 on a snow-covered path. The individual is equipped with ski poles, suggesting
345 they are on a cross-country skiing adventure. The person’s posture is relaxed, with
346 their head turned slightly to the side, as if they are taking in the beauty of their

surroundings. The forest is a dense thicket of tall, slender trees with a mix of coniferous and deciduous species. The trees are dusted with snow, their branches heavy with the white substance. The forest floor is blanketed in a thick layer of snow, untouched except for the path the skier is following. The lighting in the image is soft and diffused, indicative of an overcast sky. The colors are muted, with the red of the skier's jacket standing out against the otherwise monochromatic palette of the snowy forest. The image is rich in detail, with the texture of the snow and the bark of the trees being particularly noticeable. The composition of the image is balanced, with the skier positioned slightly off-center, following the rule of thirds. The path the skier is on leads the viewer's eye deeper into the forest, creating a sense of depth and perspective. The background is blurred, drawing focus to the skier and the immediate surroundings, but still providing a sense of the vastness of the forest. The image is a testament to the beauty of nature, captured in a moment of quiet solitude. The skier's journey through the snowy forest is a journey of exploration and connection with the natural world. The image is a masterclass in landscape photography, with its attention to detail, use of color, and composition. It is a snapshot of a moment, frozen in time, that speaks to the viewer's sense of adventure and the allure of the great outdoors. This image is a stunning example of professional photography, with its attention to detail, use of color, and composition. It is a snapshot of a moment, frozen in time, that speaks to the viewer's sense of adventure and the allure of the great outdoors. The image is a testament to the beauty of nature, captured in a moment of quiet solitude. The skier's journey through the snowy forest is a journey of exploration and connection with the natural world. The image is a masterclass in landscape photography, with its attention to detail, use of color, and composition. It is a snapshot of a moment, frozen in time, that speaks to the viewer's sense of adventure and the allure of the great outdoors. This image is a stunning example of professional photography, with its attention to detail, use of color, and composition. It is a snapshot of a moment, frozen in time, that speaks to the viewer's sense of adventure and the allure of the great outdoors. The image is a testament to the beauty of nature, captured in a moment of quiet solitude. The skier's journey through the snowy forest is a journey of exploration and connection with the natural world. The image is a masterclass in landscape photography, with its attention to detail, use of color, and composition. It is a snapshot of a moment, frozen in time, that speaks to the viewer's sense of adventure and the allure of the great outdoors. The image is a testament to the beauty of nature, captured in a moment of quiet solitude. The skier's journey through the snowy forest is a journey of exploration and connection with the natural world. The image is a masterclass in landscape photography, with its attention to detail, use of color, and composition. It is a snapshot of a moment, frozen in time, that speaks to the viewer's sense of adventure and the allure of the great outdoors. The image is a testament to the beauty of nature, captured in a moment of quiet solitude. The skier's journey through the snowy forest is a journey of exploration and connection with the natural world. The image is a masterclass in landscape photography, with its attention to detail, use of color, and composition. It is a snapshot of a moment, frozen in time, that speaks to the viewer's sense of adventure and the allure of the great outdoors.

Note: The LLaVA output exhibited significant repetition in the latter portion, repeating descriptive phrases multiple times—a common characteristic requiring downstream filtering, further supporting our claim of how increasing LLaVA output length does not lead to meaningful increase in image generation performance.

Stage 3: GPT-4.1 Prompt Optimization The lengthy LLaVA description was condensed into a focused 77-word prompt optimized for diffusion model generation:

"Lone figure in vibrant red jacket and black pants, walking away with ski poles on a snow-covered forest path, surrounded by tall coniferous and deciduous trees dusted with snow, balanced off-center composition, soft diffused lighting, muted monochromatic palette with striking red accent, professional landscape photography, serene winter atmosphere, detailed snow texture and tree bark, blurred background for depth, sense of solitude and adventure, rule of thirds, natural beauty, exploration, connection with nature, high detail."

Stage 4: Image Generation The optimized prompt was fed to SDXL, producing the output shown in Figure 7. The generated image successfully captures the core semantic elements: winter forest setting, lone figure in red jacket, snow-covered environment, and professional photography aesthetic.

B.2 Analysis

This case study demonstrates Copy2Prompt's ability to preserve semantic fidelity while avoiding direct pixel reuse. The generated image maintains the original's compositional structure and photo-



Figure 7: Generated image: SDXL output from Copy2Prompt pipeline

401 tographic style almost perfectly. However, the color pallete is slightly off as the subject is slightly
402 wearing a different color clothing the trees have different colors. We determine that this is a succesful
403 generation, as this is compositionally faithful to the original image, while producing a distinctly
404 different visual output suitable for copyright-compliant applications. Furthermore, it can be fed back
405 into similar image-to-image models with custom prompts to edit the image to be even closer to the
406 source image. For example, a sample prompt could be "most tree leaves are yellow-green, and the
407 person is wearing a black jacket with a red camping backpack."