

Leveraging Semantic Attribute Binding for Free-Lunch Color Control in Diffusion Models

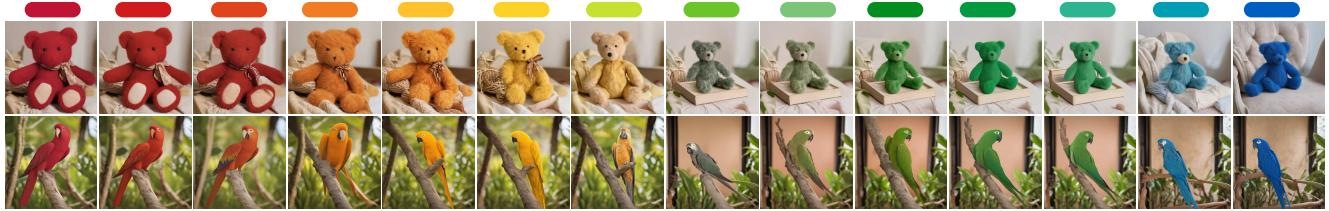
Héctor Laria^{1,2}Alexandra Gomez-Villa^{1,3}Jiang Qin⁴Muhammad Atif Butt^{1,2}Bogdan Raducanu^{1,2}Javier Vazquez-Corral^{1,2}Joost van de Weijer^{1,2}Kai Wang¹¹Computer Vision Center, Spain²Universitat Autònoma de Barcelona, Spain³Universitat de València, Spain⁴Harbin Institute of Technology, China

Figure 1. *ColorWave* accurately reproduces subtle color variations in smooth interpolation between similar tones. Each column shows a different target object rendered with gradually shifting colors (displayed above each image). The results demonstrate that our method *ColorWave* is sensitive to small changes in the RGB color space while preserving realistic object appearance and scene composition.

Abstract

Recent advances in text-to-image (T2I) diffusion models have enabled remarkable control over various attributes, yet precise color specification remains a fundamental challenge. Existing approaches, such as *ColorPeel*, rely on model personalization, requiring additional optimization and limiting flexibility in specifying arbitrary colors. In this work, we introduce *ColorWave*, a novel training-free approach that achieves exact RGB-level color control in diffusion models without fine-tuning. By systematically analyzing the cross-attention mechanisms within IP-Adapter, we uncover an implicit binding between textual color descriptors and reference image features. Leveraging this insight, our method rewrites these bindings to enforce precise color attribution while preserving the generative capabilities of pretrained models. Our approach maintains generation quality and diversity, outperforming prior methods in accuracy and applicability across diverse object categories. Through extensive evaluations, we demonstrate that *ColorWave* establishes a new paradigm for structured, color-consistent diffusion-based image synthesis.

1. Introduction

Recent advances in text-to-image (T2I) diffusion models [35, 39, 40] have revolutionized image generation, offering unprecedented control over various attributes including composition [26, 49], style [42, 45, 50], subject matter [10, 38], and visual aesthetics [12, 13]. These models can generate complex scenes with specific objects, artistic styles, and spatial arrangements based on textual descriptions. Among these controllable attributes, color stands as particularly crucial, as it fundamentally shapes viewer perception, conveys emotional tone, and often serves as a critical design specification in professional contexts from brand identity to product development [41]. Yet despite remarkable progress in generative capabilities, T2I models exhibit a significant limitation in precise color control. When users specify colors through linguistic descriptors like “red” or “blue”, these terms encompass broad ranges of potential shades, making it challenging to achieve exact color matching in generated content—a capability essential for practical design applications where color fidelity is non-negotiable.

T2I personalization methods, such as DreamBooth [38] or CustomDiffusion [21], can be used to address this lim-

itation and learn precise color prompts. However, these methods were shown to obtain unsatisfactory color generations [4], often failing to correctly disentangle shape from color information. Methods such as ColorPeel [4] enable precise color control by learning specific color prompts through additional training. However, these approaches require separate optimization processes for each individual color, resulting in significant computational overhead and limiting practicality. The need for model fine-tuning creates a fundamental barrier to flexibly specifying arbitrary colors during inference time, constraining creative applications and design tasks that demand precise color control. Furthermore, the learning-based strategy of ColorPeel becomes harder to implement in advanced T2I diffusion models, where multiple text encoders are commonly employed to process textual inputs. Models such as SDXL [33], SD3 [9], and Flux [22] introduce additional complexity, making training-based methods increasingly impractical and often ineffective for precise color control.

Recent advancements in image conditioning for diffusion models, particularly through mechanisms like IP-Adapter [45, 50, 51], have enabled new forms of control by conditioning generation on reference images. While these adapters were primarily designed for style transfer or subject-driven generation, they encode complex visual information in ways that have not been fully explored. Through systematic investigation, we discovered a previously unrecognized property: *semantic attribute binding* which refers to the implicit bindings between visual attributes and their semantic representations within multi-modal space. Specifically, we found that these bindings create a mapping between RGB values in reference images and linguistic color descriptors in text prompts—a connection that existing methods, including IP-Adapter itself, do not explicitly exploit. Our insight enables precise RGB-level color manipulation without any model fine-tuning, revealing capabilities latent within IP-Adapter that haven’t been previously identified or utilized for targeted control.

In this paper, we introduce *ColorWave*, a novel training-free approach for precise color control in diffusion models. Unlike previous methods, *ColorWave* enables users to specify exact RGB values (see Figure 1) for objects in generated images without requiring any additional training or model fine-tuning. Our method leverages semantic attribute binding capabilities within the IP-Adapter framework, effectively “rewiring” these connections to achieve precise color attribution to target objects. By introducing minimal additional modules that exploit pretrained diffusion models, we achieve unprecedented color accuracy while maintaining generation quality and diversity. To summarize, our key contributions include:

- The first training-free approach *ColorWave* for precise color control in diffusion-based generation, enabling

specification of arbitrary colors without additional optimization.

- A novel technique that exploits the binding mechanism between visual features and semantic counterparts in the latent space of diffusion models.
- A selective attention modulation strategy that preserves the generative capabilities of the base model while enabling precise color control for targeted objects.
- Comprehensive evaluations demonstrating superior performance across diverse object categories and color specifications compared to existing methods, both in terms of color accuracy and generation quality.

2. Related work

2.1. Image-Conditioned T2I Diffusion Models

Conditional image-based text-to-image (T2I) models integrate auxiliary image inputs with textual prompts, allowing for more refined control over spatial layout [7, 24], object content [6, 23], and stylistic attributes [14, 42]. For pixel-level structural guidance, ControlNet [52] employs a duplicated U-Net [36] encoder to process various input maps, while UniControlNet [54] extends this with dual adapters handling both local (pixel-level) and global (CLIP image embedding) controls, enabling simultaneous spatial and content-aware guidance.

As another technique pipeline, IP-Adapter [51] leverages the Vision Transformer (ViT) encoder from CLIP [34] to extract real image features, integrating them into the U-Net backbone via cross-attention layers for more coherent conditioning. Building on this, IPAdapter-Instruct [37] introduces an instructional prompt mechanism, enhancing interpretability and reducing the need for separate model training for different conditioning inputs, addressing a key limitation of ControlNet. However, the role of the additional cross-attention mechanism in IP-Adapter remains underexplored. In this paper, we present the first findings on its properties, revealing that IP-Adapter implicitly establishes a binding between visual attributes and their corresponding semantic representations.

2.2. Color Control in T2I diffusion models

With the rapid advancements in T2I generation, various text-guided image editing approaches [5, 15, 25, 29, 30, 47, 53] have been developed to enable controllable modifications. For instance, methods like Imagic [18] and P2P [16] leverage Stable Diffusion (SD) models for structure-preserving edits. InstructPix2Pix [3] is an extension of P2P by allowing human-like instructions for image editing. Following that, pix2pix-zero [32] propose noise regularization and *cross-attention* guidance to retrain the structure of a given image. Another technique stream which can also achieve controllable generation is transfer learning

for T2I models, or also referred to *T2I model adaptation* and *personalized generation*. It aims at adapting a given model to a *new concept* by given images from the users and bind the new concept with a unique token. As a result, the adaptation model can generate various renditions for the new concept guided by text prompts. Depending on whether the adaptation method is fine-tuning the T2I model [21, 38], or freezing the T2I backbone [8, 10, 48].

However, all these existing techniques rely heavily on the generative capacity of diffusion models and struggle to achieve fine-grained control over color attributes in image editing and generation tasks. Only a limited number of works [4, 11] have begun addressing the challenge of precise color generation. Rich-Text [11] enhances color fidelity through multi-pass processing with global and local diffusion models, but at the cost of high computational overhead and reduced color accuracy. In contrast, Color-Peel [4] introduces color prompt learning to improve color alignment with user input. However, this method requires extensive training and is restricted to handling only a few color names per training session. In this paper, we argue that the limitations of prior approaches stem from their failure to recognize the implicit attribute binding between the IP-Adapter [51] module and T2I diffusion models [33, 35]. Leveraging this property, we propose a training-free solution for precise color generation, allowing seamless adaptation to any user-specified color inputs, particularly benefiting artistic and creative applications.

3. Methodology

3.1. Preliminaries

T2I Diffusion Models. We built on the SDXL [33] model, consisting of two primary components: an autoencoder and a diffusion model ϵ_θ . The model ϵ_θ is trained by the loss:

$$L_{LDM} = \mathbb{E}_{z_0, \epsilon \sim \mathcal{N}(0,1), t} \left[\|\epsilon - \epsilon_\theta(z_t, t, \tau_\xi(\mathcal{P}))\|_2^2 \right] \quad (1)$$

where ϵ_θ is a UNet, conditioning a latent input z_t , a timestep $t \sim U(1, T)$, and a text embedding $\tau_\xi(\mathcal{P})$. More specifically, text-guided diffusion models generate an image from the textual condition as $\mathcal{C}_{text} = \tau_\xi(\mathcal{P})$, where τ_ξ is the CLIP text encoder [17]¹. The cross-attention map is derived from $\epsilon_\theta(z_t, t, \mathcal{C}_{text})$. After predicting the noise, diffusion schedulers [27, 28] are used to predict the latent z_{t-1} . As an example with the DDIM scheduler [43], the formula is:

$$z_{t-1} = \sqrt{\frac{\alpha_{t-1}}{\alpha_t}} z_t + \sqrt{\alpha_{t-1}} \left(\sqrt{\frac{1}{\alpha_{t-1}} - 1} - \sqrt{\frac{1}{\alpha_t} - 1} \right) \cdot \epsilon_\theta(z_t, t, \mathcal{C}_{text}) \quad (2)$$

where α_t is a predefined scalar function. Here we simplify the z_{t-1} inference process as $z_{t-1} = \mathcal{G}(z_t, t, \mathcal{C}_{text})$.

¹SDXL uses two text encoders and concatenate the embeddings.

IP-Adapter. Building on T2I diffusion models, the IP-Adapter [51] introduces additional controllability by conditioning the T2I model on a conditional image \mathcal{I}_{ip} . Practically, this involves leveraging a pretrained T2I diffusion model and incorporating a cross-attention layer to the (projected) image condition following each text-prompt conditioning layer. The conditional image is encoded in the low-dimensional CLIP image embedding space [17] to capture high-level semantic information. By denoting the CLIP image encoder as τ_ϕ and IP-Adapter projection as \mathbf{IP} , this process is adding a new image condition $\mathcal{C}_{img} = \mathbf{IP}(\tau_\phi(\mathcal{I}_{ip}))$ to the T2I model as $z_{t-1} = \mathcal{G}(z_t, t, \mathcal{C}_{text}, \mathcal{C}_{img})$.

Studied architecture. An overview of our approach, *ColorWave*, is provided in Figure 4. At its core, our method leverages the base architecture of a pretrained diffusion model (SDXL) with an integrated IP-Adapter. The primary insight driving our work is the discovery of *semantic attribute binding*—a phenomenon where visual attributes in reference images form implicit connections with their corresponding linguistic descriptors in text prompts. We demonstrate how this previously unexplored property can be harnessed for training-free color control in generated images.

To fully exploit this binding mechanism and achieve precise color control, *ColorWave* introduces two key enhancements to the standard IP-Adapter framework: (1) automatic color name generation, which determines the optimal linguistic color descriptor for any user-specified RGB value; and (2) spatial prior addition, which ensures accurate color attribution to target objects by refining attention maps. These components work in concert, effectively “rewiring” the cross-attention mechanism to strengthen the binding between user-specified colors and target objects while maintaining the generative capabilities of the underlying diffusion model. We detail *ColorWave* in the following sections.

3.2. Semantic attribute binding in IP-Adapters

IP-Adapter [37, 51] has demonstrated remarkable effectiveness in transferring visual attributes from reference images to generated content. In this paper, we introduce a key insight: the IP-Adapter implicitly establishes latent-space connections between visual features and their semantic counterparts. As illustrated in Figure 2, this phenomenon, which we call *semantic attribute binding*, creates a direct correspondence between colors in reference images and their linguistic descriptors in prompts. When using a single reference image containing multiple colors (Figure 2a), the model correctly extracts and applies each color based on the corresponding color name in the prompt (e.g., “green”, “yellow” or “red”). Furthermore, when the colors in the reference image are modified (Figure 2b), the generated results reflect these changes while maintaining the binding to the appropriate color names. This binding even works with synthetic color patches (Figure 2c), allowing for

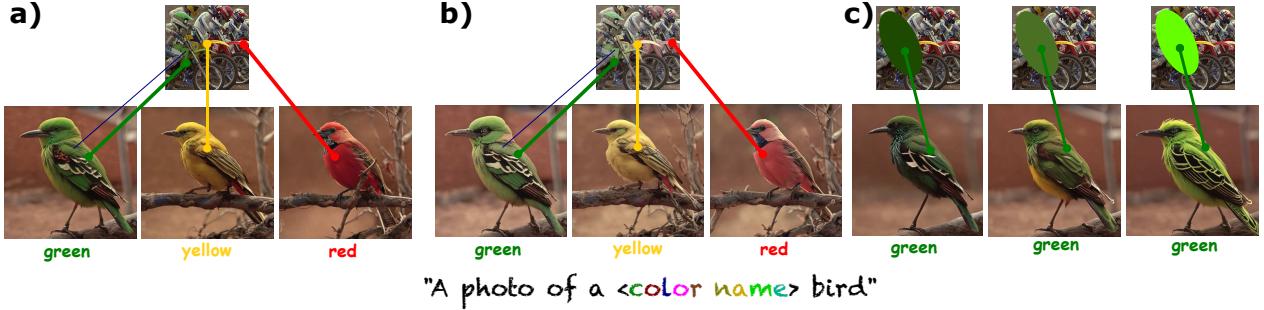


Figure 2. **Illustration of semantic attribute binding.** On top, the color-guidance image is provided. a) Given the color name used in the prompt, the generated results will pick a different color from the color-guidance image. b) Changing the colors of the color guidance image, results in similar changes in the generated images. c) The exact color which should be generated for the used color name can also be a synthetic example.

precise RGB-level control. These observations reveal that the adapter inherently associates specific color values with their linguistic descriptors, enabling exact color specification without requiring any additional training or fine-tuning.

Color name attribution. *Semantic attribute binding* stems from the architectural design of IP-Adapter, particularly its decoupled cross-attention mechanism. As shown by Ye et al. [51], the IP-Adapter introduces separate cross-attention layers for image features while keeping the original text cross-attention layers intact. Crucially, both cross-attention mechanisms share the same query matrices, while maintaining separate key and value projection matrices for text and image features:

$$\mathbf{x} = \text{softmax} \left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d}} \right) \mathbf{V} + \text{softmax} \left(\frac{\mathbf{Q}\mathbf{K}'^\top}{\sqrt{d}} \right) \mathbf{V}'. \quad (3)$$

This sharing of query projection matrices creates an implicit alignment between the text and image feature spaces. During IP-Adapter training, the model learns to map visual color information from the image encoder to the same latent space that the text encoder uses to represent color concepts. Specifically, when a colored image is provided as input, the IP-Adapter learns to bind the visual color properties with their corresponding linguistic color representations that the diffusion model already understands.

To effectively leverage this semantic attribute binding, we need to determine which color names in natural language correspond to specific RGB values. Following foundational work in color linguistics by Berlin and Kay [2], we recognize that human languages typically use a limited set of basic color terms to describe the entire color spectrum. In English, there are eleven basic color terms, with eight being chromatic (red, green, blue, yellow, brown, orange, pink, and purple).

To validate our hypothesis that the key projections may be collapsing different color modalities into the same representation, we compute the similarity matrix between color

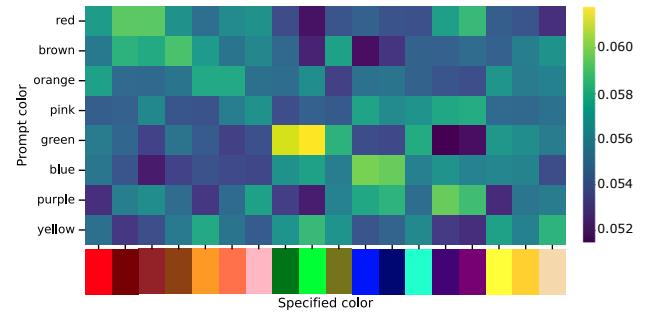


Figure 3. **Semantic attribute binding.** Phenomenon visualization through a similarity matrix between color names in text prompts and RGB color values. The heatmap shows the normalized dot product similarity between key projections of color word tokens and image features.

names used in text prompts and RGB color values in reference images. We calculate the projection inner product $\langle \mathbf{K}, \mathbf{K}' \rangle$ to examine similarities between color word tokens from text prompts and RGB feature tokens from the adapter. This similarity measure captures the strength of the semantic binding between textual color descriptors and visual color representations.

Figure 3 visualizes this similarity matrix for various color names and specified RGB values. The heatmap reveals clear patterns of correspondence, with notably higher similarity values along the diagonal where color names match their expected RGB values (e.g., the “green” token shows strongest binding with green RGB values). This confirms our hypothesis that the IP-Adapter implicitly establishes these bindings during its training process, even though it was never explicitly trained for color matching.

Interestingly, we observe that the similarity values are not uniform across all color pairs. Some color terms (like “green” and “blue”) show stronger and more distinctive binding patterns, while others exhibit more diffuse relationships. This variation aligns with findings in color linguistics that certain basic color terms have more centralized and

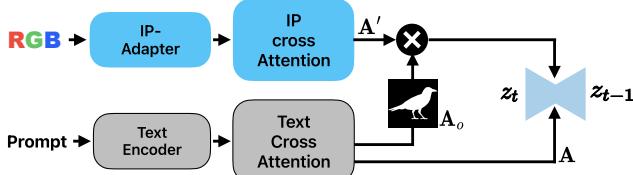


Figure 4. **Overview of ColorWave**. Our approach leverages semantic attribute binding between IP-Adapter and text cross-attention pathways to achieve precise color control. User-specified RGB values are encoded through IP-Adapter and selectively bound to object tokens in the text prompt, enabling training-free color attribution while preserving generative quality.

consistent representations across languages and visual systems [19].

By quantifying these semantic bindings, we can automatically determine the appropriate color name to use in prompts based on user-specified RGB values, creating a bidirectional mapping between the continuous color space and discrete linguistic color categories. This mapping forms the foundation of our color control approach, allowing us to exploit the implicit color knowledge already encoded in the model’s attention mechanisms.

Direct Semantic Attribute Binding Limitations. While the binding mechanism offers a promising avenue for color control, naive implementations face several challenges, as illustrated in Figure 5:

- **Shape and Size Sensitivity:** As shown in Figure 5b,c, changing the shape or size of the color reference region alters the resulting object colors. Different image statistics from these variations influence the generation process, leading to inconsistent color attribution.
- **Ambiguous Color Attribution:** When multiple regions contain the desired color (Figure 5d), the binding mechanism struggles to correctly attribute the target color, often selecting one region at random. Furthermore, all colors present in the reference image influence the final appearance of the generated object.
- **Context Limitation:** IP-Adapter generates new images based on the statistics of the reference image, restricting the diversity of backgrounds and contexts (Figure 5e). To generate an object with a specific color in various contexts, each reference image would need to contain that exact color, severely limiting creative freedom.
- **Synthetic Reference Limitation:** Using purely synthetic color references (Figure 5e) results in flat, unrealistic surfaces lacking natural variations, highlights, and textures, producing results that appear artificial and fail to capture the nuanced appearance of real-world objects.

To address these challenges, our *ColorWave* approach introduces a selective attention modulation strategy that pre-

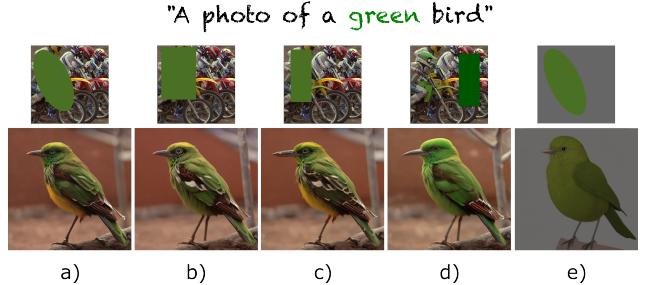


Figure 5. **Limitations when directly exploiting semantic attribute binding for color control.** (a) Reference image with oval. (b,c) Shape and size variations alter the resulting bird coloration despite using the same green color. (d) With multiple green regions, attribution becomes ambiguous and inconsistent. (e) Using synthetic color references produces flat, unrealistic textures.

serves the generative capabilities of the base model while enabling precise color control for targeted objects.

3.3. Spatial Prior Addition

In our approach *ColorWave*, we leverage the decoupled cross-attention design established in previous adapters [45, 51] for its minimal interference with pretrained knowledge and compatibility with other methods. We implement a strategy of minimal yet precise intervention in the generative model to achieve effective color binding. First, we inject the user-specified color information at a single optimal point in the network, where it most effectively influences the desired subject, following successful approaches in prior work [45, 46]. Second, we rely on the model’s inherent understanding to determine proper color placement based on its pretrained knowledge.

The challenge is that color information lacks spatial specificity, so we need to ensure it is applied to the correct object. To solve this, we query the model for the target object’s location and constrain our intervention specifically to that region, ensuring the specified color binds precisely to the relevant semantic token in the prompt. As illustrated in Fig. 4, we implement this by masking the adapter’s contribution \mathbf{A}' with the attention map of the desired object \mathbf{A}_o before incorporating it into the final features \mathbf{x} :

$$\mathbf{x} = \mathbf{A} + [\mathbf{A}_o] \mathbf{A}' \quad (4)$$

Here, $\mathbf{A} \in \mathcal{R}^{h \times H \times W}$ represents the original cross-attention output, while the operator $[\cdot]$ averages the attention maps across attention heads and preserves only the highest-value regions. This selective approach complements and ensures effective semantic attribute binding described earlier between the specified color and the target object, while minimizing unintended color attribution to other elements in the scene. The evolution of the object mask throughout the generation process is visualized in the appendix.

4. Experiments

4.1. Experimental setup

Dataset. There is no standardized dataset specifically designed for evaluating precise color control in diffusion models. For fair and direct comparison with prior work, we adopt the evaluation framework established by ColorPeel [4], which defines two color generation tasks:

- **Coarse-grained color set:** Four basic colors (red, green, blue, and yellow) with predefined RGB values.
- **Fine-grained color set:** Eighteen more specialized colors including "salmon", "beige", "navy", and "indigo", each with specific RGB values.

For both sets, we use identical prompts to those used in ColorPeel, including objects like "a [color] bowl on the table", "a [color] teddy bear in Times Square", and "a [color] sofa in living room". Each prompt is rendered with 20 different random seeds to account for generation variability, resulting in a comprehensive evaluation suite of 200 images for the coarse-grained set and 360 images for the fine-grained set.

It is important to note that while previous methods like ColorPeel require separate training procedures for each new target color, our *ColorWave* approach inherently accepts any arbitrary RGB value without additional training. The distinction between "coarse" and "fine" grained colors is only maintained here for comparative purposes — our method treats all colors equally, handling any point in the RGB color space with the same level of precision. To demonstrate this capability, we include additional experiments with randomly selected colors outside both predefined sets, showing that *ColorWave* maintains consistent performance across the entire color spectrum without any per-color optimization.

Evaluation metrics. We adopt the evaluation metrics proposed by ColorPeel [4]. Specifically, we compute several complementary metrics:

- **Euclidean Distance in CIE Lab color space** (Δ_E and $\Delta_{E_{Ch}}$ when luminance is removed): These metrics measure perceptual uniformity between generated and target colors; lower values indicate better color matching.
- **Mean Angular Error (MAE) in sRGB:** quantifies color deviation in terms of chromaticity, helping to understand differences in hue and saturation independent of intensity.
- **Mean Angular Error (MAE) in Hue:** This analyzes the difference between target and generated colors irrespective of brightness and saturation.

For metric and implementation details refer to appendix.

Comparison methods. Our method *ColorWave* represents the first training-free approach for precise color control

in diffusion models, establishing a distinct methodological category from existing approaches that require color-specific training. Our primary comparison is with ColorPeel [4], the current state-of-the-art in color-specific prompt learning. For completeness, we compare against all baselines evaluated in ColorPeel, that includes *training-free* T2I generation baselines: (1) vanilla Stable Diffusion with color name text prompts; (2) Rich-Text [11], which enhances adherence to complex text descriptions; And also *training-based* personalization methods besides ColorPeel [4]: (3) Textual Inversion [10], learning new pseudo-words in embedding space; (4) DreamBooth [38], fine-tuning the entire diffusion model; and (5) Custom Diffusion [21], optimizing projection matrices in cross-attention layers. The training-free nature of *ColorWave* creates a fundamental asymmetry in this comparison – while ColorPeel and other training-based approaches require separate optimization processes for each individual color or complex training regimes to handle multiple colors simultaneously, *ColorWave* inherently accepts any arbitrary RGB triplet without modification. This represents not just an incremental improvement but a paradigm shift in how precise color control can be achieved in generative models, transcending the limitations of discrete color vocabularies and model-specific optimizations.

4.2. Qualitative Comparisons

To demonstrate the versatility and effectiveness of *ColorWave*, we present qualitative results across several key capabilities: precise color attribution across diverse objects, fine-grained color control, and generalization to complex color patterns and textures. These experiments highlight the flexibility of our training-free approach in scenarios that would require extensive additional training in previous methods. More results can be found in the Supplementary.

Arbitrary color attribution. We evaluate *ColorWave*'s ability to precisely apply user-specified colors to a wide range of subjects including inanimate objects, animals, plants, and human clothing. Figure 6 demonstrates that our method successfully generates images where target objects accurately reflect the desired reference colors. Notably, *ColorWave* maintains high visual quality across diverse scenarios, preserving natural lighting effects, surface properties, and contextual integration while achieving precise color matching.

Fine color navigation. Beyond coarse color attribution, we examine whether our method can detect and reproduce subtle color variations. Figure 1 showcases results from smooth interpolation between similar color tones. This precision demonstrates the sensitivity of semantic attribute binding and shows that the proposed method can navigate

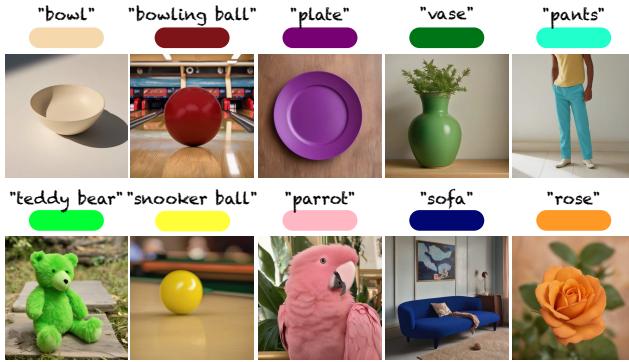


Figure 6. **Arbitrary color attribution across diverse subjects.** *ColorWave* precisely applies user-specified colors to various target objects while maintaining natural lighting, material properties, and contextual integration.



Figure 7. **Generalizability to complex color patterns and textures.** These examples illustrate how semantic attribute binding extends beyond simple color matching to more sophisticated visual attribute control.

the continuous color space with remarkable accuracy.

ColorWave generalizability. We further investigate the extensibility of our approach beyond single-color attribution. Figure 7 demonstrates two advanced applications: (1) applying complete color palettes to target objects, and (2) transferring complex textures while maintaining precise color control. In both cases, the method preserves the structural integrity and material properties of the target objects while successfully integrating the specified color patterns.

4.3. Quantitative Analysis

We compare *ColorWave* with training-based and training-free baselines. Such quantitative results are demonstrated in Table 1. For each method, we generated images, and extract the mask of the object, following the same pipeline as in ColorPeel [4]. Percentages in MAE metrics denote the percentage of pixels inside the mask used for the computation (selecting those closest to the ground truth). This table clearly shows the superiority of *ColorWave*. Particularly, *ColorWave* achieved notably lower ΔE error in CIE Lab color space as compared to the existing training-free methods and it is only worse than the training-based state-of-the-art ColorPeel, which indicates that *ColorWave* generates perceptually better colors than its direct competitors. In addition, *ColorWave* also achieved comparatively much lower mean angular error in both sRGB and Hue when compared to the other training-free methods, which signifies a higher degree of color accuracy in terms of chromaticity and hue in our generated images.

User study We conducted a user study with 15 participants to perceptually evaluate our results, comparing *ColorWave* against ColorPeel [4], TI [10], Rich-Text [11], DB [38], and CD [21]. We followed the same experimental paradigm as the one defined in ColorPeel.

The experiment took place in a controlled lab environment to ensure reliability. All participants were tested for correct color vision using the Ishihara test. The study followed a two-alternative forced choice (2AFC) method. Observers viewed three images on a monitor set to RGB: the central image represented the target color, while the left and right images displayed results generated by our method and one of the competing methods, with their positions randomized. We tested the same 10 different prompts and four colors (red, green, blue, and yellow) defined in ColorPeel.

To analyze the results, we compared *ColorWave* against each competing method using the Thurstone Case V Law of Comparative Judgment model [44]. This approach provided us with z-scores and a 95% confidence interval, calculated using the method proposed in [31]. The results, shown in Fig. 8, indicate that *ColorWave* is statistically significantly better than all five competing algorithms, including ColorPeel, which learns a specific token for each color. These findings highlight *ColorWave*'s effectiveness in generating more realistic and accurate colors given an RGB triplet.

4.4. Ablations

Color name attribution. We investigate the constraints of semantic attribute binding by intentionally creating discrepancies between the color mentioned in the prompt and the color specified as input to our model. Figure 9 illustrates these results using a consistent prompt “*A photo of a red parrot*” while varying the input color values. Our findings

Table 1. Quantitative comparison with baselines over various evaluation metrics. All numbers are the smaller the better (\downarrow). The best results of both *training-based* and *training-free* technique streams are highlighted in bold. Training time is provided in the last column.

Method	ΔE	ΔE_{ch}	MAE (sRGB)			MAE (Hue)			Time (min)
			10%	50%	100%	10%	50%	100%	
Training Based	TI [10]	48.98	44.29	15.22	19.51	23.90	52.66	69.35	90.88
	DB [38]	50.71	46.29	14.75	19.30	23.70	47.12	67.13	88.72
	CD [21]	48.47	42.23	13.43	17.93	22.43	31.63	55.07	78.43
	ColorPeel (3D) [4]	21.39	16.51	4.36	7.76	12.08	2.63	6.47	21.35
	ColorPeel (2D) [4]	20.45	15.29	4.83	7.88	12.13	3.18	7.43	21.46
Training Free	SD [35]	47.45	41.55	12.89	20.04	26.93	30.17	54.14	86.38
	Rich-Text [11]	36.62	32.48	9.91	13.29	18.53	50.55	72.77	93.51
	ColorWave (Ours)	29.91	26.57	8.65	10.39	12.71	13.69	22.14	33.12

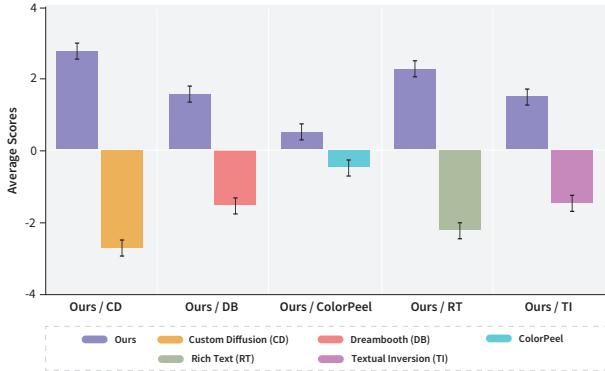


Figure 8. **User study.** Z-scores are higher for better human preference. Our method outperforms baselines and the state-of-the-art color control generation method ColorPeel [4].

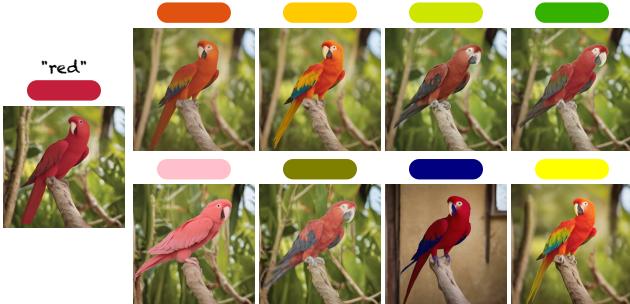


Figure 9. **Boundaries of semantic attribute binding when prompt and target colors diverge.** All images were generated using the prompt “A photo of a red parrot” while varying the input color. This demonstrates that semantic attribute binding operates within color neighborhoods, with effectiveness decreasing as perceptual distance increases between prompt color and target color.

reveal a semantic proximity effect, understood by Fig. 3; the success of color attribution depends on the perceptual and semantic distance between the prompt color and the target color. As shown in the upper row of Figure 9, when the

target color (orange) shares perceptual similarity with the prompt color (red), our method successfully binds the new color to the parrot. However, when the target color (green) represents an opponent color on the color wheel, the binding fails, and the model defaults to generating a parrot closer to the prompted red color. The lower row demonstrates this effect across a broader range of colors. Pink, which shares perceptual properties with red (both are warm colors with similar hue angles), successfully binds to the parrot. In contrast, semantically distant colors such as olive, navy, and yellow fail to override the prompted “red” descriptor, resulting in parrots that retain reddish tones despite the different color input.

4.5. Limitations

Despite the effectiveness of *ColorWave* in achieving training-free color control, several limitations remain. First, our approach relies on the quality of cross-attention maps generated by the underlying diffusion model. When these maps exhibit attention leakage—where attention for specific objects spreads to unintended regions—our method’s color attribution precision diminishes. While state-of-the-art diffusion models continue to improve attention mechanism quality, certain complex scenes or ambiguous prompts may still produce suboptimal attribution maps.

Furthermore, our current implementation faces challenges when attributing different colors to multiple objects in a single scene. The cross-attention maps tend to degrade in quality as the number of attribution targets increases, potentially leading to color bleeding between objects. This limitation becomes particularly evident in dense scenes with numerous objects requiring distinct color specifications.

For future works, we intend to address these limitations by exploring robust attention map generation techniques and developing specialized approaches for multi-object color attribution.

5. Conclusion

We introduced *ColorWave*, the first training-free approach for precise color control in text-to-image (T2I) diffusion models. By leveraging the latent binding between visual attributes and semantic representations within the IP-Adapter, we enabled direct RGB-level color manipulation without requiring additional model fine-tuning. Our method effectively rewrites cross-attention mechanisms to establish precise color attribution while maintaining the generative quality and diversity of pretrained diffusion models. Extensive evaluations demonstrated that our method significantly outperforms prior training-free approaches in color accuracy and object-specific fidelity. Our method also shows competitive performance versus state-of-the-art, fine-tuning methods like ColorPeel, thus overcoming the limitations of prompt-based color descriptions which are computationally expensive. Our findings reveal previously untapped capabilities in existing diffusion models, paving the way for new directions in controllable generation.

Acknowledgements

We acknowledge project PID2022-143257NB-I00, financed by MCIN/AEI/10.13039/501100011033 and FEDER, and the Generalitat de Catalunya CERCA Program. JVC was funded by Grant PID2021-128178OB-I00 funded by MCIN/AEI/10.13039/501100011033, ERDF “A way of making Europe” and the Departament de Recerca i Universitats from Generalitat de Catalunya with ref. 2021SGR01499. This work was partially supported by the grant Càtedra ENIA UAB-Cruïlla (TSI-100929-2023-2) from the Ministry of Economic Affairs and Digital Transition of Spain.

References

- [1] Aishwarya Agarwal, Srikrishna Karanam, Tripti Shukla, and Balaji Vasan Srinivasan. An image is worth multiple words: Multi-attribute inversion for constrained text-to-image synthesis. *ICML*, 2024. 11
- [2] Brent Berlin and Paul Kay. *Basic color terms: Their universality and evolution*. Univ of California Press, 1991. 4
- [3] Tim Brooks, Aleksander Holynski, and Alexei A Efros. Instructpix2pix: Learning to follow image editing instructions. In *CVPR*, 2023. 2
- [4] Muhammad Atif Butt, Kai Wang, Javier Vazquez-Corral, and Joost van de Weijer. Colorpeel: Color prompt learning with diffusion models via color and shape disentanglement. In *ECCV*, 2024. 2, 3, 6, 7, 8
- [5] Songyan Chen and Jiancheng Huang. Fec: Three finetuning-free methods to enhance consistency for real image editing. In *International Conference on Image Processing, Computer Vision and Machine Learning*, 2023. 2
- [6] Wenhui Chen, Hexiang Hu, Yandong Li, Nataniel Rui, Xuhui Jia, Ming-Wei Chang, and William W Cohen. Subject-driven text-to-image generation via apprenticeship learning. In *NeurIPS*, 2023. 2
- [7] Hongsuk Choi, Isaac Kasahara, Selim Engin, Moritz Graule, Nikhil Chavan-Dafle, and Volkan Isler. FineControlNet: Fine-level Text Control for Image Generation with Spatially Aligned Text Control Injection. *arxiv.org/abs/2312.09252*, 2023. 2
- [8] Ziyi Dong, Pengxu Wei, and Liang Lin. Dreamartist: Towards controllable one-shot text-to-image generation via contrastive prompt-tuning. *arXiv preprint arXiv:2211.11337*, 2022. 3
- [9] 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 *ICML*, 2024. 2
- [10] Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or. An image is worth one word: Personalizing text-to-image generation using textual inversion. *ICLR*, 2023. 1, 3, 6, 7, 8
- [11] Songwei Ge, Taesung Park, Jun-Yan Zhu, and Jia-Bin Huang. Expressive text-to-image generation with rich text. In *ICCV*, 2023. 3, 6, 7, 8
- [12] Daniel Geng, Inbum Park, and Andrew Owens. Visual anagrams: Generating multi-view optical illusions with diffusion models. In *CVPR*, 2024. 1
- [13] Alexandra Gomez-Villa, Kai Wang, Alejandro C Parraga, Bartłomiej Twardowski, Jesus Malo, Javier Vazquez-Corral, and Joost van de Weijer. The art of deception: Color visual illusions and diffusion models. *CVPR*, 2025. 1
- [14] Zhen Han, Chaojie Mao, Zeyinzi Jiang, Yulin Pan, and Jingfeng Zhang. Stylebooth: Image style editing with multimodal instruction. *arXiv preprint arXiv:2404.12154*, 2024. 2
- [15] Amir Hertz, Kfir Aberman, and Daniel Cohen-Or. Delta denoising score. In *ICCV*, 2023. 2
- [16] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Prompt-to-prompt image editing with cross attention control. *ICLR*, 2023. 2
- [17] Gabriel Ilharco, Mitchell Wortsman, Nicholas Carlini, Rohan Taori, Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hannaneh Hajishirzi, Ali Farhadi, and Ludwig Schmidt. Open clip, 2021. 3, 11
- [18] Bahjat Kawar, Shiran Zada, Oran Lang, Omer Tov, Huiwen Chang, Tali Dekel, Inbar Mosseri, and Michal Irani. Imagic: Text-based real image editing with diffusion models. *CVPR*, 2023. 2
- [19] Paul Kay and Terry Regier. Language, thought and color: recent developments. *Trends in cognitive sciences*, 10(2): 51–54, 2006. 5
- [20] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *ICCV*, 2023. 11
- [21] Nupur Kumari, Bingliang Zhang, Richard Zhang, Eli Shechtman, and Jun-Yan Zhu. Multi-concept customization of text-to-image diffusion. *CVPR*, 2023. 1, 3, 6, 7, 8

- [22] Black Forest Labs. Flux. <https://github.com/black-forest-labs/flux>, 2024. 2
- [23] Dongxu Li, Junnan Li, and Steven Hoi. Blip-diffusion: Pre-trained subject representation for controllable text-to-image generation and editing. In *NeurIPS*, 2024. 2
- [24] Ming Li, Taojiannan Yang, Huafeng Kuang, Jie Wu, Zhaoning Wang, Xuefeng Xiao, and Chen Chen. ControlNet++: Improving Conditional Controls with Efficient Consistency Feedback. *ECCV*, 2024. 2
- [25] Senmao Li, Joost van de Weijer, Taihang Hu, Fahad Shahbaz Khan, Qibin Hou, Yaxing Wang, and Jian Yang. Styledifusion: Prompt-embedding inversion for text-based editing. *arXiv preprint arXiv:2303.15649*, 2023. 2
- [26] Nan Liu, Shuang Li, Yilun Du, Antonio Torralba, and Joshua B Tenenbaum. Compositional visual generation with composable diffusion models. In *ECCV*, 2022. 1
- [27] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. DPM-Solver++: Fast Solver for Guided Sampling of Diffusion Probabilistic Models. *NeurIPS*, 2022. 3
- [28] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. DPM-Solver: A Fast ODE Solver for Diffusion Probabilistic Model Sampling in Around 10 Steps. In *NeurIPS*, 2022. 3
- [29] Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jianjun Wu, Jun-Yan Zhu, and Stefano Ermon. SDEdit: Guided image synthesis and editing with stochastic differential equations. In *ICLR*, 2022. 2
- [30] Ron Mokady, Amir Hertz, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Null-text inversion for editing real images using guided diffusion models. *CVPR*, 2023. 2
- [31] Ethan D Montag. Empirical formula for creating error bars for the method of paired comparison. *J. Elec. Imag.*, 15(1): 010502–010502, 2006. 7
- [32] Gaurav Parmar, Krishna Kumar Singh, Richard Zhang, Yijun Li, Jingwan Lu, and Jun-Yan Zhu. Zero-shot image-to-image translation. *Proceedings of the ACM SIGGRAPH Conference on Computer Graphics*, 2023. 2
- [33] 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. In *ICLR*, 2024. 2, 3, 11
- [34] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, 2021. 2
- [35] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022. 1, 3, 8
- [36] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *MICCAI*, 2015. 2
- [37] Ciara Rowles, Shimon Vainer, Dante De Nigris, Slava Elizarov, Konstantin Kutsy, and Simon Donné. Ipadapter-instruct: Resolving ambiguity in image-based conditioning using instruct prompts, 2024. 2, 3
- [38] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *CVPR*, 2023. 1, 3, 6, 7, 8
- [39] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S Sara Mahdavi, Rapha Gontijo Lopes, et al. Photorealistic text-to-image diffusion models with deep language understanding. *NeurIPS*, 2022. 1
- [40] Alex Shonenkov, Misha Konstantinov, Daria Bakshandaeva, Christoph Schuhmann, Ksenia Ivanova, and Nadiia Klokova. Deepfloyd-if. <https://github.com/deep-floyd/IF>, 2023. 1
- [41] Satyendra Singh. Impact of color on marketing. *Management decision*, 44(6):783–789, 2006. 1
- [42] Kihyuk Sohn, Nataniel Ruiz, Kimin Lee, Daniel Castro Chin, Irina Blok, Huiwen Chang, Jarred Barber, Lu Jiang, Glenn Entis, Yuanzhen Li, et al. Styledrop: Text-to-image generation in any style. *NeurIPS*, 2023. 1, 2
- [43] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising Diffusion Implicit Models. *ICLR*, 2021. 3
- [44] Louis L Thurstone. A law of comparative judgment. In *Scaling*, pages 81–92. Routledge, 1927. 7
- [45] Haofan Wang, Matteo Spinelli, Qixun Wang, Xu Bai, Zekui Qin, and Anthony Chen. Instantstyle: Free lunch towards style-preserving in text-to-image generation. *arXiv preprint arXiv:2404.02733*, 2024. 1, 2, 5, 11
- [46] Haofan Wang, Peng Xing, Renyuan Huang, Hao Ai, Qixun Wang, and Xu Bai. Instantstyle-plus: Style transfer with content-preserving in text-to-image generation. *arXiv preprint arXiv:2407.00788*, 2024. 5, 11
- [47] Kai Wang, Fei Yang, Shiqi Yang, Muhammad Atif Butt, and Joost van de Weijer. Dynamic prompt learning: Addressing cross-attention leakage for text-based image editing. *NeurIPS*, 2023. 2
- [48] Kai Wang, Fei Yang, Bogdan Raducanu, and Joost van de Weijer. Multi-class textual-inversion secretly yields a semantic-agnostic classifier. In *WACV*, 2025. 3
- [49] Ruichen Wang, Zekang Chen, Chen Chen, Jian Ma, Haonan Lu, and Xiaodong Lin. Compositional text-to-image synthesis with attention map control of diffusion models. In *AAAI*, 2024. 1
- [50] Peng Xing, Haofan Wang, Yanpeng Sun, Qixun Wang, Xu Bai, Hao Ai, Renyuan Huang, and Zechao Li. Csgo: Content-style composition in text-to-image generation. *arXiv 2408.16766*, 2024. 1, 2
- [51] Hu Ye, Jun Zhang, Sibo Liu, Xiao Han, and Wei Yang. Ipadapter: Text compatible image prompt adapter for text-to-image diffusion models. *arXiv preprint arXiv:2308.06721*, 2023. 2, 3, 4, 5, 11
- [52] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *ICCV*, pages 3836–3847, 2023. 2
- [53] Shiwen Zhang, Shuai Xiao, and Weilin Huang. Forgedit: Text guided image editing via learning and forgetting. *arXiv preprint arXiv:2309.10556*, 2023. 2

- [54] Shihao Zhao, Dongdong Chen, Yen-Chun Chen, Jianmin Bao, Shaozhe Hao, Lu Yuan, and Kwan-Yee K Wong. Uni-controlnet: All-in-one control to text-to-image diffusion models. In *NeurIPS*, 2023. 2

A. Color binding regions as a function of steps

Figure 10 displays color binding regions based on $\lceil \mathbf{A}_o \rceil$ every 10-th denoising step,



Figure 10. Color binding regions based on $\lceil \mathbf{A}_o \rceil$. Depicted every 10-th denoising step.

B. Experimental details

Evaluation metric details For each generated image, we use the Segment-Anything model [20] to automatically extract masks for the target objects, allowing precise measurement of color attributes only within relevant regions. We report these metrics under multiple conditions: considering all pixels within the object mask (100%), and also considering the top 10% and 50% of pixels closest to the target color, which helps account for natural variations in object appearance such as highlights and shadows.

Implementation details We implement *ColorWave* using a pretrained Stable Diffusion XL (SDXL) [33] model as our base architecture. For the image conditioning component, we utilize the IP-Adapter framework [51] with an encoder based on OpenCLIP-ViT-H-14 [17]. We only inject the color embedding into *the first decoder layer* to compute cross-attention maps, which shows better *stylization* performance as previous works proved [1, 45, 46].

For the adapter masking, we keep the top 20% largest values on each object map. During inference, we process the user-specified RGB values that creates temporary color reference images. These references are encoded by the IP-Adapter and strategically injected into the model’s cross-attention layers.

C. Extra results of *ColorWave*

Extra results are showcased in Figures 11 to 18.



Figure 11. "A photo of a red -object-" - Reference color is depicted in the external frame of the grid - Objects same as in ColorPeel method (see qualitative evaluation)



Figure 12. "A photo of a red -object-" - Reference color is depicted in the external frame of the grid - Objects same as in ColorPeel method (see qualitative evaluation)



Figure 13. "A photo of a pink -object-" - Reference color is depicted in the external frame of the grid - Objects same as in ColorPeel method (see qualitative evaluation)



Figure 14. "A photo of a orange -object-" - Reference color is depicted in the external frame of the grid - Objects same as in ColorPeel method (see qualitative evaluation)



Figure 15. "A photo of a purple -object-" - Reference color is depicted in the external frame of the grid - Objects same as in ColorPeel method (see qualitative evaluation)



Figure 16. "A photo of a green -object-" - Reference color is depicted in the external frame of the grid - Objects same as in ColorPeel method (see qualitative evaluation)



Figure 17. "A photo of a blue -object-" - Reference color is depicted in the external frame of the grid - Objects same as in ColorPeel method (see qualitative evaluation)



Figure 18. "A photo of a blue -object-" - Reference color is depicted in the external frame of the grid - Objects same as in ColorPeel method (see qualitative evaluation)