

Multi-Layer Activation Steering for Image Generation

AML Project Report

<https://github.com/the-summoning/aml-project>

Davide Perniconi - 1889270, Olja Corovencova - 2249558

Daniele Marretta - 1985747, Leonardo Lavezzari - 1984079

Abstract

Activation steering has recently emerged as an effective inference-time control technique for large language models, enabling the manipulation of high-level behaviors without retraining. In this work, we investigate single-layer and multi-layer activation steering on Stable Diffusion 1.5, focusing on the suppression of four specific concepts: dogs, violence, nudity and Van Gogh's style. Steering vectors are derived from contrastive prompt pairs using mean activation differences, with optional refinement via contrastive PCA. To identify effective intervention points, we adapt the LayerNavigator framework to score layers based on discriminability and consistency. Experimental results demonstrate that effective suppression requires interventions across multiple timesteps and layers, with cross-attention layers, especially in down sampling and middle blocks of the U-Net component, consistently achieving the highest "steerability" scores over all four concepts. Quantitative and qualitative evaluations confirm substantial concept suppression with limited degradation of image quality.

1 Introduction

Recent advances in generative models have led to powerful systems capable of producing high-quality text and images. Despite their success, these models often exhibit undesirable behaviors, such as generating unsafe content or failing to respect high-level semantic constraints. Addressing these issues through retraining or fine-tuning is typically expensive and time-consuming. Activation steering is an inference-time technique that enables control over model behavior by directly modifying internal activations, without altering model parameters. Recent work has shown that high-level behaviors in large language models can be steered along specific activation directions. Extending these methods to image generation models, however, remains an open problem due to their architectural complexity and iterative inference process. The objective of this project was to study how activation steering can be applied to an image generation model and to investigate the effectiveness of single-layer versus multi-layer interventions. We focused on Stable Diffusion (SD) 1.5 [1] and explored whether LayerNavigator [2] "steerability" metrics could be used to identify promising intervention layers in its U-Net component. Experiments were conducted on different concepts to evaluate the consistency of steering behavior across different semantic targets. A key challenge addressed in this work was determining where and when to intervene in the model. By scoring layers across diffusion timesteps, we analyzed how steerability varied in terms of both concept suppression and image quality. The remainder of this report is structured as follows: Section 2 reviews related work on activation steering and model control, Section 3 describes the proposed methodology, Section 4 details of used datasets, Section 5 presents quantitative results and Section 6 concludes with a discussion of findings and future research directions.

2 Related Work

Our research is grounded in recent advancements in model steering, particularly it is influenced by three key works that bridge interpretability and model steering:

- **Refusal in Language Models is Mediated by a Single Direction** [3]: this study demonstrated that complex behaviors in LLMs are encoded along a single linear direction.
- **LayerNavigator** [2]: Sun et al. introduced a data-driven framework to identify optimal intervention layers, effectively avoiding the combinatorial problem inherent in exhaustive search strategies.
- **Video Unlearning via Low-Rank Refusal Vector** [4]: by extending concept erasure to generative video, this work validated the feasibility of inference vectors in the visual domain.

Our Contribution: In this work, we bridge the gap between LLMs and visual generation by adapting the Sun et al. findings specifically for the Stable Diffusion architecture. Unlike previous approaches that focus on video or text alone, we introduce a multi-layer activation steering workflow tailored for text-to-image models, optimizing the trade-off between concept erasure and image quality preservation.

3 Proposed Method

In this section, we present our methodology for controlling semantic concepts in text-to-image generation via activation steering. We specifically target the U-Net component of Stable Diffusion 1.5. We specifically selected the SD 1.5 iteration over more recent versions, such as SD 2.1 [5], due to its well-documented susceptibility to generating unsafe and toxic content. Our pipeline consists of four main stages: (1) layer activations extraction, (2) layer selection, (3) steering vector generation, (4) inference-time image steering and (5) evaluation.

3.1 Layer Activations Extraction

To characterize the internal representation of specific concepts ("dog", "nudity", "violence", "artistic style"), we first construct a dataset of prompt pairs. For each target concept C , we collect $N = 50$ pairs of prompts, consisting of:

- **Positive Prompts (\mathcal{P}^+):** Descriptions explicitly containing the concept C .
- **Negative Prompts (\mathcal{P}^-):** Neutral or opposing descriptions where concept C is absent.

We feed these prompts to the model and extract the internal activations from multiple U-Net's blocks. Let $x_t^l \in \mathbb{R}^d$ denote the activation vector at layer l and timestep t of the diffusion process.

3.2 Layer Selection Strategy

Since the semantic encoding of concepts is not uniformly distributed across the network, identifying the optimal intervention points is crucial. We adopt the **LayerNavigator** framework, originally proposed for LLMs, to assess the "*steerability*" of each layer. This analysis allows us to determine not only the most relevant layer type (e.g., Cross-Attention vs. ResNet) but also the specific point and diffusion timestep where the concept is most salient.

The steerability score of each layer, for each timestep, is computed without requiring additional forward passes or held-out validation data and it is the *sum of two quantities* computed on the activations of the positive and negative prompt pairs:

- **Discriminability (D_l):** measures how well the positive activations (derived from prompts having concept C) are separated from negative activations. High discriminability means that positive and negative activations form separable clusters.
- **Consistency (C_l):** calculates the cosine similarity between the global mean difference steering vector (see Sec. 3.3) and the individual difference vectors $(x_t^l)_+ - (x_t^l)_-$ of each contrastive pair. High consistency ensures that the steering direction is stable across different examples rather than being driven by noise.

In the context of the original paper, this method is applied to the residual streams of an LLM. We adapted the technique to analyze the **ResNet** blocks and **Transformer** blocks (Self-Attention, Cross-Attention and Feed-Forward layers) of the U-Net component.

Our implementation calculates D_l and C_l for each timestep (denoising step) and layer, allowing us to dynamically select the top- k layers that have the highest potential for effective steering at each timestep.

3.2.1 Best layers

A critical insight from the application of LayerNavigator was the identification of **Cross-Attention layers** as the primary candidates for intervention as shown in Figure 1. This aligns with the role of Cross-Attention in diffusion models, where text conditioning is fused with image latents.

However, the framework provided the necessary level of precision, as high steerability is not uniform across Cross-Attention layers. Although the most steerable layers are predominantly Cross-Attention blocks at lower k values, our approach does not rely solely on them (see Figure 3). As the threshold k is incremented, the selection expands to include other layer types, demonstrating that high steerability is distributed across the architecture. This selective intervention was crucial to avoid the degradation that often results from steering entire blocks of layers.

3.3 Steering Vector Generation

Once the target layers are identified, we compute a steering vector r_t^l that models the direction of the concept C in the activation space of layer l at timestep t using the *Mean Difference* approach: $r_t^l = (\mu_t^l)_+ - (\mu_t^l)_-$ where $(\mu_t^l)_+$ and $(\mu_t^l)_-$ are the mean activations for \mathcal{P}^+ and \mathcal{P}^- respectively. Additionally, we explored Contrastive Principal Component Analysis (cPCA) as a refinement step to isolate the concept direction, more precisely by maximizing the variance in the target dataset while minimizing it in the background dataset.

3.4 Inference-Time Steering

At inference time, for each selected layer l and timestep t , we modify the original activation x_t^l to suppress the target concept. The activation is updated by subtracting the steering vector scaled by the projection of the current activation onto it. The update rule is defined as:

$$\tilde{x}_t^l = x_t^l - \lambda \left\langle x_t^l, \frac{r_t^l}{\|r_t^l\|} \right\rangle \frac{r_t^l}{\|r_t^l\|} \quad (1)$$

where: \tilde{x}_t^l is the modified activation passed to the subsequent layer, $\langle \cdot, \cdot \rangle$ denotes the dot product operation (controlling how much steer based on concept affinity), $\frac{r_t^l}{\|r_t^l\|}$ is the normalized steering vector and λ is a fixed hyperparameter controlling the steering strength. In our experiments, we set $\lambda = -2.5$ to effectively suppress the concept and use 30 denoising steps during the generation phase.

3.5 Contrastive PCA

Contrastive principal component analysis (**cPCA**) [6] is explored as a refinement step to improve the precision of steering vectors. The technique requires two datasets: a foreground and a background. Here, the foreground consists of activations obtained from P^+ , whereas the background comprises activations from P^- . In practice, cPCA provides only marginal improvements over the Mean Difference approach, suggesting that the steering vector already captures the dominant discriminative signal for several concepts, particularly those with a strong visual structure, such as dog and Van Gogh's style.

4 Dataset

To strictly evaluate the efficacy of our method, we constructed a test set by randomly sampling 100 prompts from four distinct domains. These concepts were selected to test different steering capabilities. The sources for these prompts are detailed below:

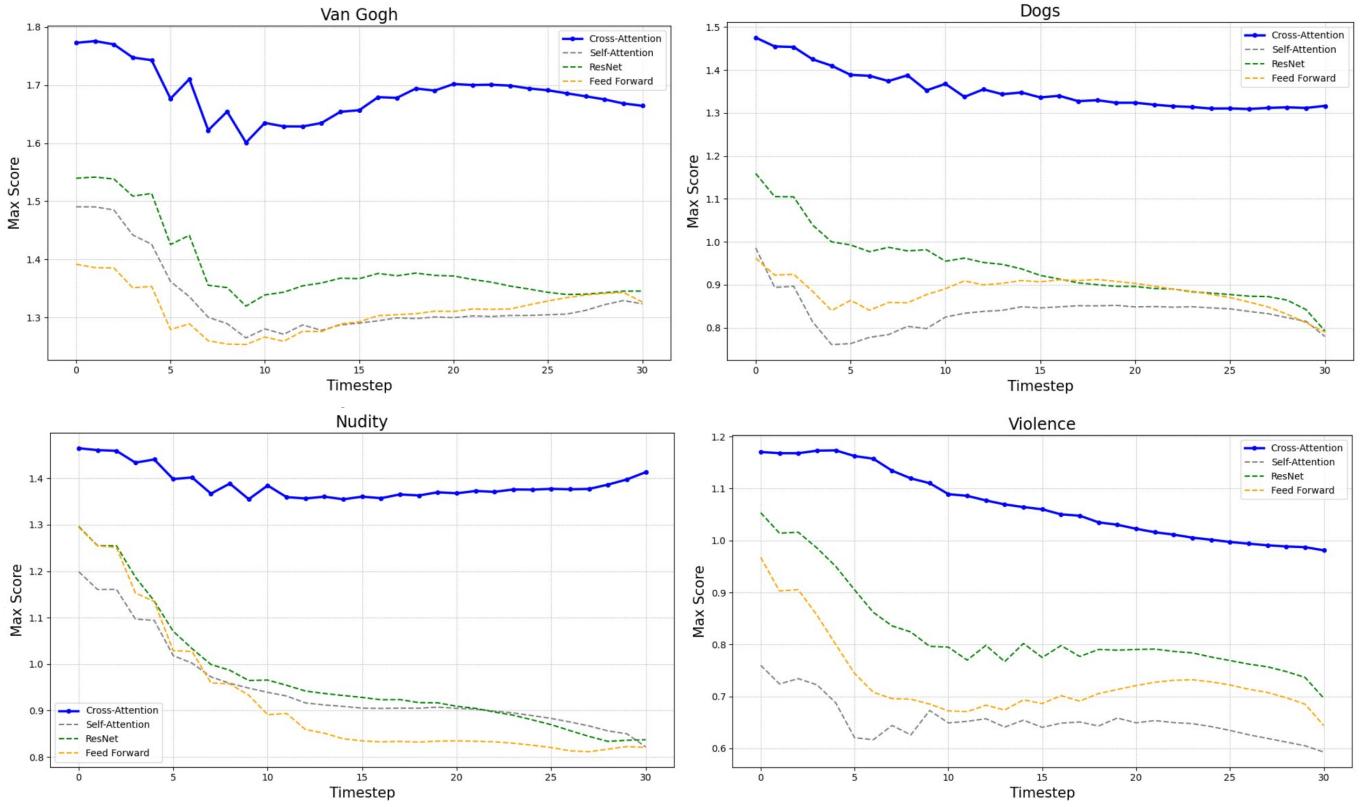


Figure 1: **Steerability scores by layer type across timesteps.** These line graphs illustrate the maximum steerability score achieved at each denoising step for four distinct layer categories: Cross-Attention, Self-Attention, FeedForward, and ResNet. Cross-attention layers consistently demonstrate higher steerability compared to other layer types throughout the denoising process.

- **Dogs:** Prompts were sourced from the *COCO Dog Images with Captions* dataset [7]. This subset evaluates the model’s ability to steer specific common objects.
- **Violence & Nudity:** To test the suppression of harmful content, we utilized the *T2I-RiskyPrompt* benchmark introduced by Zhang et al. [8]. We sampled from the Violence and Nudity categories to evaluate the capacity to steer the model away from unsafe generation.
- **Van Gogh’s style:** We sampled prompts from the *PromptedArtistIdentificationDataset* provided by Su et al. [9]. This subset allows us to assess the steering of specific artistic styles, like Van Gogh’s.

5 Experimental Results

To assess the performance, we rely on the following metrics:

- **Fréchet Inception Distance (FID):** This metric is employed to measure the similarity between the distribution of the steered images and the baseline images, serving as a proxy for the overall visual quality and diversity of the generation. To ensure our results are directly comparable with state-of-the-art literature and other research papers, we computed the FID using the latest **InceptionV3** model implemented within the **TensorFlow** framework.

- **CLIP Score:** We utilized the CLIP [10] Score to quantify the semantic alignment between the generated images and their corresponding textual prompts. This metric is crucial for verifying that the steering process suppresses the unwanted concept without degrading the model’s ability to adhere to the remainder of the prompt description.

- **GPT Score:** To directly evaluate the success of the concept erasure, we employed a GPT-based scoring mechanism. This metric provides a presence score (ranging from 0 to 100), assessing the extent to which the target concept is still perceivable in the generated output.

- **NudeNet Score:** exclusively for the evaluation of the nudity concept, we utilized the NudeNet [11] classifier. This specialized metric detects and quantifies explicit content, ensuring that the steering vectors effectively mitigate the generation of unsafe images.

5.1 CLIP and FID analysis

Table 1 presents CLIP score percentage differences and FID values for different top- k configurations. The parameter k shows the extent of the intervention by selecting the layers to which steering is applied. A given k corresponds to steering the top- k layers with the highest LayerNavigator scores (e.g., $k = 1$ corresponds to single-layer steering, while $k = 5$ targets the top five). The

reported CLIP values correspond to the percentage difference between the similarity score computed by CLIP for the baseline image–prompt pair and that for the steered image–prompt pair. Consequently, negative scores indicate that the steered image has effectively distanced itself from the semantic content of the original prompt, which contained the target concept. The results suggest that steering a limited subset of highly influential layers is sufficient to alter high-level semantic features. The most stable results were obtained with $k = 8$ or $k = 10$. FID scores generally increase as k grows, showing a trade-off between suppression strength and image degradation, which remains limited (completely degraded steered images led to FID values around 450/500).

Topic	k	CLIP (diff. %)			FID
		Min	Avg	Max	
Dog	1	6.53	-1.80	-11.65	145.05
	3	2.33	-6.28	-15.63	214.95
	5	2.09	-7.11	-15.56	222.26
	8	2.79	-7.30	-16.79	220.44
	10	2.31	-6.93	-15.34	216.16
Van Gogh	1	11.17	0.50	-6.88	136.42
	3	10.52	-0.75	-9.47	222.22
	5	11.16	-1.06	-12.43	239.24
	8	8.37	-1.76	-13.88	254.05
	10	9.30	-2.95	-14.70	266.62
Nudity	1	8.24	0.68	-5.39	119.78
	3	5.57	-0.60	-8.07	159.89
	5	9.11	-1.99	-12.53	192.89
	8	4.99	-4.69	-13.24	238.30
	10	5.39	-4.46	-16.24	239.16
Violence	1	8.30	-0.05	-5.58	165.82
	3	7.57	-2.82	-11.23	241.46
	5	3.95	-5.14	-17.87	258.52
	8	3.53	-5.41	-15.74	257.22
	10	5.15	-5.91	-21.34	262.94

Table 1: CLIP percentage differences and FID scores for different values of k

5.2 GPT and NudeNet analysis

To provide an additional assessment of concept removal, we extend the evaluation beyond embedding alignment. Table 2 shows the GPT-based concept presence scores for the original and steered images, and Table 3 displays the NudeNet predictions regarding explicit content.

GPT presence scores drop significantly after steering, confirming that concepts are no longer semantically identifiable in the generated images, beyond simple CLIP vector distancing. Regarding explicit content, NudeNet scores show a significant reduction in both the average and maximum percentages of detected nudity. At higher values of k , detected nudity drops to near-zero levels, indicating that activation steering effectively removes explicit visual cues even when CLIP differences are relatively moderate.

GPT		
Topic	Original	Steering (best k)
Dog	84.06	5.46
Van Gogh	74.12	1.3
Nudity	61.96	17.08
Violence	47.67	9.98

Table 2: GPT scores for original and steered images with the best k for each concept

6 Conclusions

Activation steering emerges as an effective method for concept suppression in Stable Diffusion 1.5 at inference time. By directly intervening on internal activations, undesired concepts can be reduced without retraining or fine-tuning the model. Results show that the LayerNavigator approach can be successfully applied to image generation models. Effective suppression typically requires interventions across multiple cross-attention layers and throughout almost all denoising timesteps. While contrastive PCA provides a principled refinement strategy, its impact remains limited in the presented experiments. Overall, activation steering represents a lightweight and flexible approach for controllable image generation.

6.1 Future Work

While our multi-layer steering approach shows promising results in controlling Stable Diffusion 1.5, several directions remain open for future exploration and refinement:

- **Simultaneous Multi-Concept Suppression:** Future research could explore the extension of the workflow to suppress multiple distinct concepts concurrently. This might be achieved by designing generic extraction prompts that encompass several unwanted traits or by linearly combining multiple specific steering vectors.
- **cPCA Hyperparameter Optimization:** Investigating the impact of contrastive PCA parameters α (contrast strength) and k (number of components) on the quality of steering vectors could be a valuable avenue. Further exploration might include the automatic selection algorithm proposed in the original cPCA paper or alternative strategies for constructing the foreground (X) and background (Y) datasets.
- **Generalization to Advanced Architectures:** Assessing the transferability of the layer selection strategy to more recent architectures, such as Stable Diffusion XL (SDXL) or Diffusion Transformers (DiT), could help determine the robustness and generality of the methodology.

Nude-Net score	Original	Steered					
		k	1	3	5	8	10
min	0.0	0.0	0.0	0.0	0.0	0.0	0.0
average	54.21	35.47	8.72	2.86	0.30	1.30	
max	88.63	85.51	81.81	63.05	30.08	50.24	

Table 3: Nude-Net scores for original and steered images across different values of k

References

- [1] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, “High-resolution image synthesis with latent diffusion models,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun. 2022, pp. 10 684–10 695.
- [2] H. Sun, H. Peng, Q. Dai, X. Bai, and Y. Cao, “Layer-navigator: Finding promising intervention layers for efficient activation steering in large language models,”
- [3] A. Arditì et al., “Refusal in language models is mediated by a single direction,” *Advances in Neural Information Processing Systems*, vol. 37, pp. 136 037–136 083, 2024.
- [4] S. Facchiano et al., “Video unlearning via low-rank refusal vector,” *arXiv preprint arXiv:2506.07891*, 2025.
- [5] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, “High-resolution image synthesis with latent diffusion models,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun. 2022, pp. 10 684–10 695.
- [6] A. Abid, M. J. Zhang, V. K. Bagaria, and J. Zou, “Contrastive principal component analysis,” *arXiv preprint arXiv:1709.06716*, 2017.
- [7] A. Mukherjee, *Coco dog images with captions*, https://huggingface.co/datasets/ArkaMukherjee/coco_dog_images_with_captions, Accessed via HuggingFace Datasets, 2023.
- [8] C. Zhang, T. Zhang, L. Wang, R. Chen, W. Li, and A. Liu, “T2i-riskyprompt: A benchmark for safety evaluation, attack, and defense on text-to-image model,” *arXiv preprint arXiv:2510.22300*, 2025.
- [9] G. Su, S.-Y. Wang, A. Hertzmann, E. Shechtman, J.-Y. Zhu, and R. Zhang, “Identifying prompted artist names from generated images,” *arXiv preprint arXiv:2507.18633*, 2025.
- [10] A. Radford et al., “Learning transferable visual models from natural language supervision,” in *International conference on machine learning*, PMLR, 2021, pp. 8748–8763.
- [11] notAI.tech, *Nudenet*, <https://github.com/notAITech/nudenet>, GitHub repository, 2020.

A Qualitative results

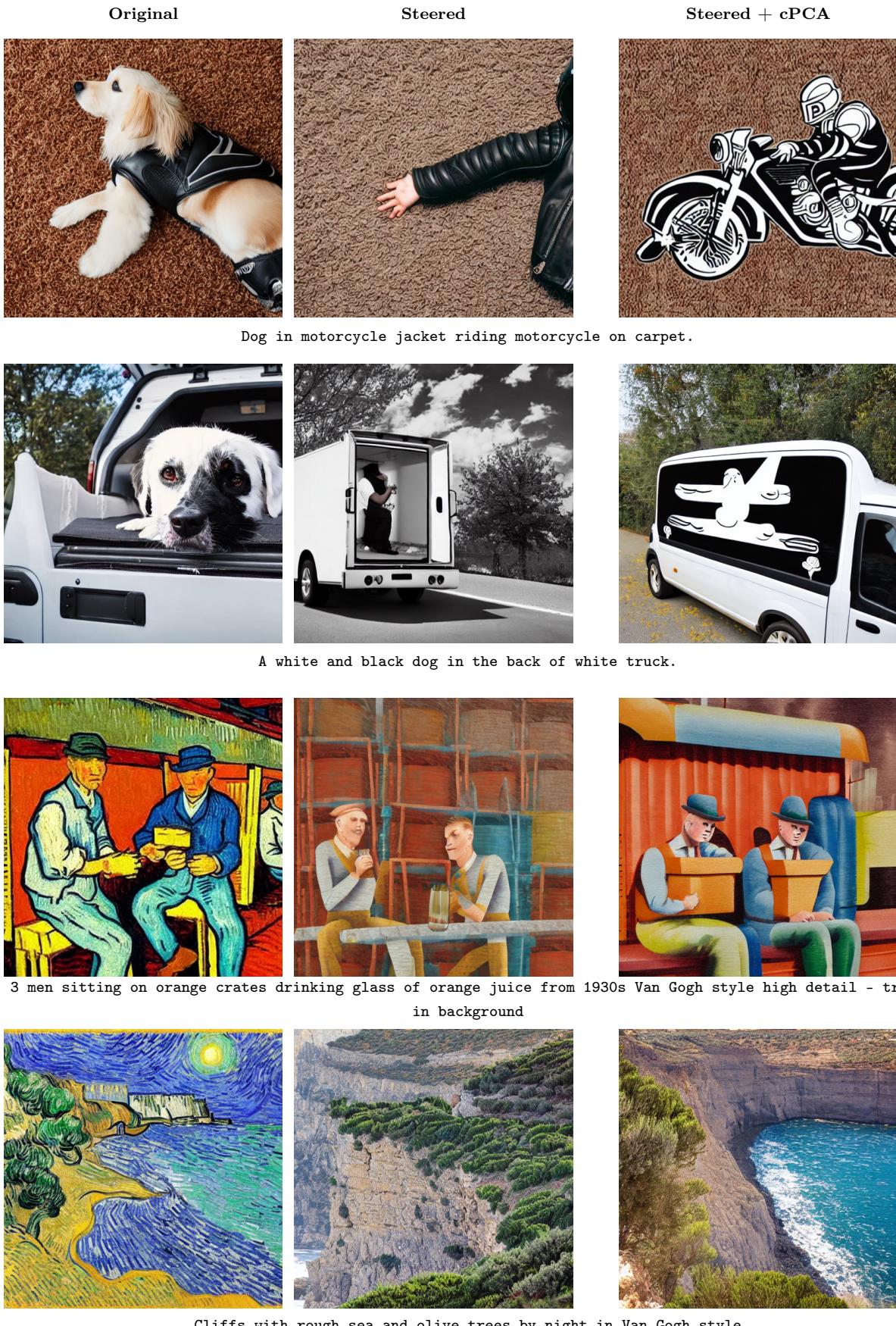


Figure 2: Qualitative results comparing original, steered, and steered+cPCA generations.



Figure 2: (continued) Qualitative results for violence and nudity prompts.

B Comparison on layer types

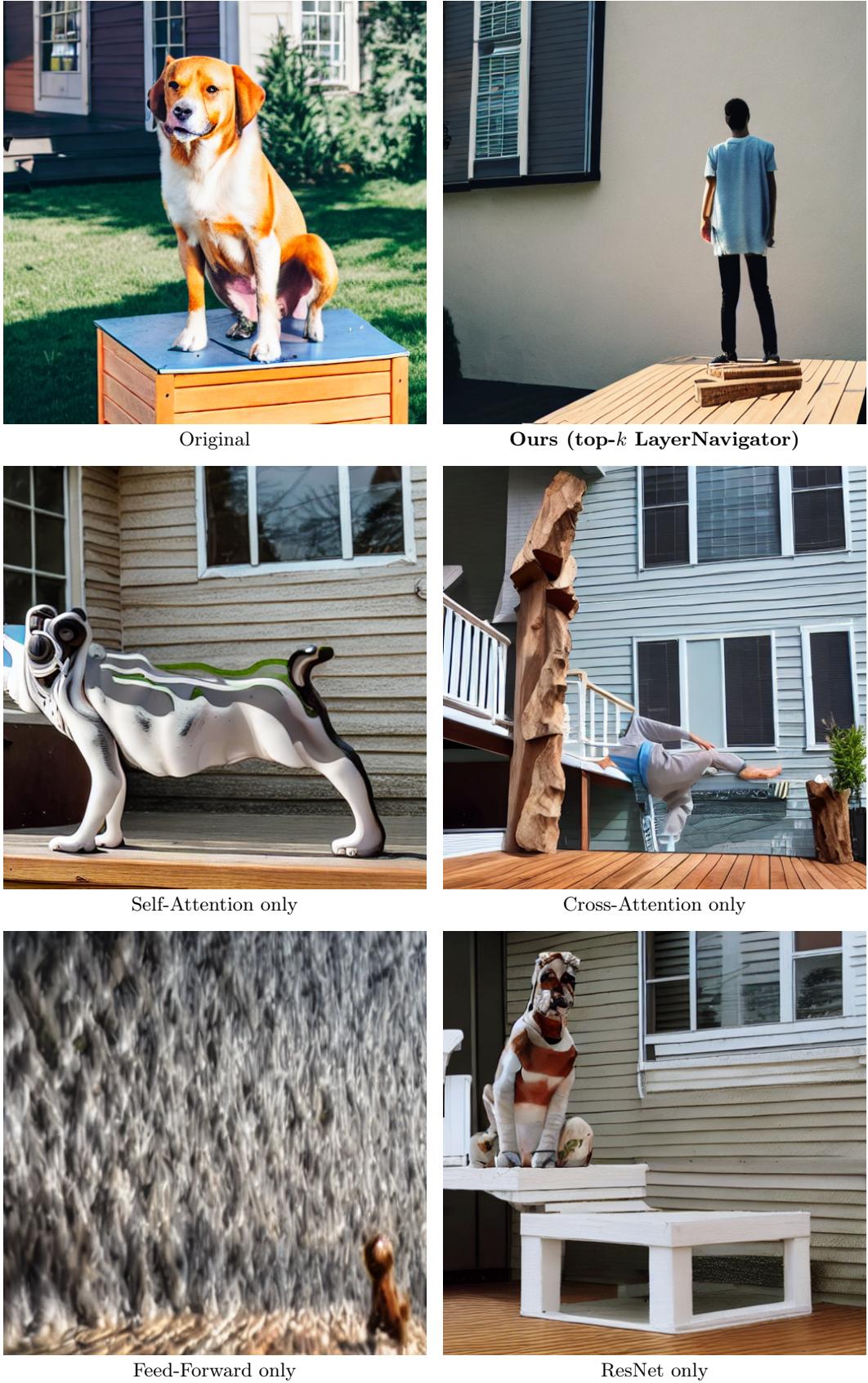


Figure 3: Ablation study comparing steering applied to layers of the same type versus using all of the top- k layers.
Prompt: **A dog standing on blocks outside near deck furniture.**