

Unpacking SDXL Turbo: Interpreting Text-to-Image Models with Sparse Autoencoders

Viacheslav Surkov Chris Wendler Mikhail Terekhov Justin Deschenaux
Robert West Caglar Gulcehre
EPFL



Figure 1. Enabling features learned by sparse autoencoders results in interpretable changes in SDXL Turbo’s image generation process. The image captions correspond to feature codes comprised of transformer block name and feature index.

Abstract

Sparse autoencoders (SAEs) have become a core ingredient in the reverse engineering of large-language models (LLMs). For LLMs, they have been shown to decompose intermediate representations that often are not interpretable directly into sparse sums of interpretable features, facilitating better control and subsequent analysis. However, similar analyses and approaches have been lacking for text-to-image models. We investigated the possibility of using SAEs to learn interpretable features for SDXL Turbo, a few-step text-to-image diffusion model. To this end, we train SAEs on the updates performed by transformer blocks within SDXL Turbo’s denoising U-net. We find that their learned features are interpretable, causally influence the generation process, and reveal specialization among the blocks. In particular, we find one block that deals mainly with image composition, one that is mainly responsible for adding local details, and one for color, illumination, and style. Therefore, our case study on SDXL Turbo is an important first step towards better understanding the internals of generative text-to-image models and showcases the potential of features learned by SAEs for the visual domain.

We provide code to reproduce our experiments, along with a demo application, in <https://github.com/surkovv/sd-xl-unbox>¹

¹The code base contains scripts to collect training data and train SAEs. Additionally, the provided demo application that allows to browse pre-trained SAE features and to edit the image generation process by modifying SAE feature coefficients.

1. Introduction

Text-to-image generation is a rapidly evolving field. The DALL-E model first captured public interest [41], combining learned visual vocabularies with sequence modeling to produce high-quality images based on user input prompts. Today’s best text-to-image models are largely based on text-conditioned diffusion models [3, 37, 38, 44–46]. This can be partially attributed to the stable training dynamics of diffusion models, which makes them easier to scale than previous approaches such as generative adversarial neural networks [16]. As a result, they can be trained on internet scale image-text datasets like LAION-5B [47] and learn to generate photorealistic images from text.

However, the underlying logic of the neural networks that enable the text-to-image pipelines we have today, due to their black-box nature, is not well understood. Unfortunately, this lack of interpretability is typical in the deep learning field. For example, advances in image recognition [26] and language modeling [6, 15] come mainly from scaling models [21], rather than from an improved understanding of their internals. Recently, the emerging field of mechanistic interpretability has sought to alleviate this limitation by reverse engineering visual models [34] and transformer-based LLMs [40]. At the same time, diffusion models have remained underexplored.

This work focuses on SDXL Turbo, a recent open-source few-step text-to-image diffusion model. We import methods from a toolbox originally developed for language models, which allows inspection of the intermediate results of the forward pass [5, 8, 12, 20]. Moreover, some of these

methods even enable reverse engineering of the entire task-specific subnets [32]. In particular, *sparse autoencoders* (SAEs) [5, 12, 54] are considered a breakthrough in interpretability for LLMs. They have been shown to decompose intermediate representations of the LLM forward pass – often difficult to interpret due to *polysemanticity*² – into sparse sums of interpretable and monosemantic features. These features are learned in an unsupervised way, can be automatically annotated using LLMs [7], and facilitate subsequent analysis, for example, circuit extraction [32].

Contributions. In this work, we investigate whether we can use SAEs to draw insights about the computation performed by the one-step generation process of SDXL Turbo, which is a recent open-source few-step text-to-image diffusion model.

To facilitate our analysis, we developed a library called *SDLens* that allows us to cache and manipulate intermediate results of SDXL Turbo’s forward pass. We use our library to create a dataset of SDXL Turbo’s intermediate feature maps of several transformer blocks inside SDXL Turbo’s U-net on 1.5M LAION-COCO prompts [47, 48]. We then use these feature maps to train multiple SAEs for each transformer block. Finally, we perform a quantitative and qualitative analysis of the SAEs’ learned features:

1. We empirically show the potential of SAEs to learn highly interpretable features that causally affect the generation in SDXL Turbo (see Fig. 1).
2. We develop visualization techniques to analyze the interpretability and causal effects of the learned features.
3. We perform a case study in which we visualize and interpret the active features in different transformer blocks, finding evidence that certain transformer blocks of SDXL Turbo’s U-net specialize in *image composition*, *adding details*, and *style*.³
4. We follow up our qualitative case study by designing multiple quantitative experiments showing that our hypotheses also hold up on larger sample sizes.
5. As part of our quantitative analysis, we create an automatic feature annotation pipeline for the transformer block, which appears responsible for *image compositions*.

Thus, we show that SAEs learn interpretable features that causally affect SDXL Turbo’s image generation process. Importantly, the learned features provide insight into the computational details of SDXL Turbo’s forward pass, such as the different roles of the investigated transformer blocks. By open-sourcing our library and SAEs, we lay the foundation for further research in this area.

²A phenomenon where a single neuron or feature encodes multiple, unrelated concepts [17].

³The blocks `down.2.1` and `up.0.1` have been already known in the community as “composition” and “style” blocks [50]. However, in this paper we provide the first thorough and fine-grained investigation of them.

Note on visualizations. Qualitative analysis by inspection of generated images is crucial for this type of research. Thus, we provide additional visualizations in the App. B and C.

2. Background

2.1. Sparse Autoencoders

Let $h(x) \in \mathbb{R}^d$ be an intermediate result during a forward pass of a neural network on input x . In a fully connected neural network, $h(x)$ could correspond to a vector of neuron activations. In transformers, which are neural network architectures that combine attention with fully connected layers and residual connections, $h(x)$ could either refer to the content of the residual stream after a layer, an update to the residual stream by a layer, or a vector of neuron activations within a fully connected layer.

It has been shown [5, 12, 54] that in many neural networks, especially LLMs, intermediate representations can be well approximated by sparse sums of $n_f \in \mathbb{N}$ learned feature vectors, i.e.,

$$h(x) \approx \sum_{\rho=1}^{n_f} s_{\rho}(x) \mathbf{f}_{\rho}, \quad (1)$$

where $s_{\rho}(x)$ are the input-dependent coefficients, most of which are equal to zero and $\mathbf{f}_1, \dots, \mathbf{f}_{n_f} \in \mathbb{R}^d$ is a learned dictionary of feature vectors. Importantly, the features are usually interpretable.

Sparse autoencoders. To implement the sparse decomposition from equation 1, the vector s containing the n_f coefficients of the sparse sum, is parameterized by a single linear layer followed by an activation function, called the *encoder*,

$$s = \text{ENC}(h) = \sigma(W^{\text{ENC}}(h - b_{\text{pre}}) + b_{\text{act}}), \quad (2)$$

in which $h \in \mathbb{R}^d$ is the latent that we aim to decompose, $\sigma(\cdot)$ is an activation function, $W^{\text{ENC}} \in \mathbb{R}^{n_f \times d}$ is a learnable weight matrix and b_{pre} and b_{act} are learnable bias terms. We omitted the dependencies $h = h(x)$ and $s = s(h)$, which are clear from the context.

Similarly, the learnable features are parametrized by a single linear layer called *decoder*,

$$h' = \text{DEC}(s) = W^{\text{DEC}} s + b_{\text{pre}}, \quad (3)$$

in which $W^{\text{DEC}} = (\mathbf{f}_1 | \dots | \mathbf{f}_{n_f}) \in \mathbb{R}^{d \times n_f}$ is a learnable matrix. Its columns take the role of learnable features and b_{pre} is a learnable bias term.⁴

⁴An extended version of this section, including training details, is in the App. E.

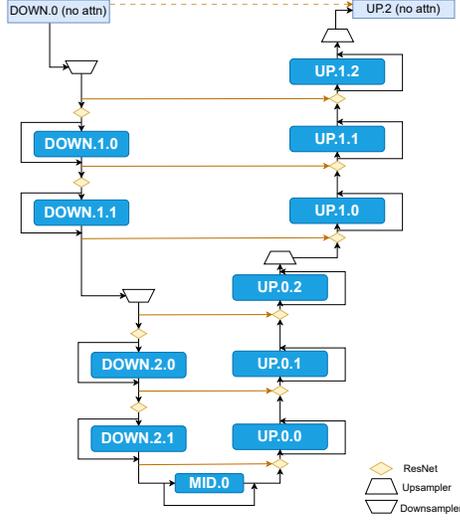


Figure 2. Cross-attention transformer blocks in SDXL’s U-net.

2.2. Few-Step Diffusion Models: SDXL Turbo

SDXL Turbo [46] is a distilled version of Stable Diffusion XL [38], a powerful latent diffusion model. SDXL Turbo allows high-quality sampling in as few as 1-4 steps. It employs a denoising network implemented using a U-net similar to Rombach et al. [44].

The U-net consists of a down-sampling path, a bottleneck, and an up-sampling path, each comprising one or more U-net blocks connected via up- and down-samplers. U-net blocks are built from residual networks, with some blocks incorporating multiple cross-attention transformer blocks while others do not. We refer to these transformer blocks by their short names (e.g., down.2.1, as illustrated in Fig. 2). Each transformer block is composed of several basic transformer layers, each consists of self-attention, cross-attention, and MLP layers. Importantly, the text conditioning is achieved via cross-attention to text embeddings performed by 11 transformer blocks embedded in the down-, up-sampling paths, and the bottleneck. An architecture diagram displaying the relevant blocks along with their names can be found in Fig. 2.

3. Sparse Autoencoders for SDXL Turbo

With the necessary definitions at hand, in this section we show a way to apply SAEs to SDXL Turbo. In the following, we assume that all SDXL Turbo generations are done using a 1-step process.

Where to apply the SAEs. We apply SAEs to updates performed within the cross-attention transformer blocks responsible for incorporating the text prompt (for details consider the App. A). Each of these blocks consists of multi-

ple transformer layers, which attend to all spatial locations (self-attention) and to the text prompt embeddings (cross-attention).

Formally, the ℓ th cross-attention transformer block updates its inputs in the following way

$$D[\ell]_{ij}^{out} = D[\ell]_{ij}^{in} + \mathcal{T}[\ell](D[\ell]_{ij}^{in}, c)_{ij}, \quad (4)$$

in which $D[\ell]_{ij}^{in}, D[\ell]_{ij}^{out} \in \mathbb{R}^{h \times w \times d}$ denote the residual stream before and after application of the ℓ -th cross-attention transformer block respectively. The transformer block itself calculates the function $\mathcal{T}[\ell] : \mathbb{R}^{h \times w \times d} \rightarrow \mathbb{R}^{h \times w \times d}$. Note that we omitted the dependence on input noise z_t and text embedding c for both $D[\ell]_{ij}^{in}(z_t, c)$ and $D[\ell]_{ij}^{out}(z_t, c)$.

We train SAEs on the residual updates $\mathcal{T}[\ell](D[\ell]_{ij}^{in}, c)_{ij} \in \mathbb{R}^d$ denoted by

$$\Delta D[\ell]_{ij} := \mathcal{T}[\ell](D[\ell]_{ij}^{in}, c)_{ij} = D[\ell]_{ij}^{out} - D[\ell]_{ij}^{in}. \quad (5)$$

That is, we jointly train one encoder $\text{ENC}[\ell]$ and decoder $\text{DEC}[\ell]$ pair per transformer block ℓ and share it over all spatial locations i, j . For notational convenience we omit block indices from now. We do this for the 4 (out of 11) transformer blocks (Fig. 2) that we found have the highest impact on the generation, namely, down.2.1, mid.0, up.0.0 and up.0.1.

Feature maps. We refer to $\Delta D \in \mathbb{R}^{h \times w \times d}$ as dense feature map and applying ENC to all image locations results in the *sparse feature map* $S \in \mathbb{R}^{h \times w \times n_f}$ with entries

$$S_{ij} = \text{ENC}(\Delta D_{ij}). \quad (6)$$

We refer to the feature map of the ρ th learned feature using $S^\rho \in \mathbb{R}^{h \times w}$. This feature map S^ρ contains the spatial activations of the ρ th learned feature. Its associated feature vector $\mathbf{f}_\rho \in \mathbb{R}^d$ is a column in the decoder matrix $W^{\text{DEC}} = (\mathbf{f}_1 | \dots | \mathbf{f}_{n_f}) \in \mathbb{R}^{d \times n_f}$. Using this notation, we can represent each element of the dense feature map as a sparse sum

$$\Delta D_{ij} \approx \sum_{\rho=1}^{n_f} S_{ij}^\rho \mathbf{f}_\rho, \text{ with } S_{ij}^\rho = 0 \text{ for most } \rho \in \{1, \dots, n_f\}. \quad (7)$$

Training. In order to train an SAE for a transformer block, we collected dense feature maps ΔD_{ij} from SDXL Turbo one-step generations on 1.5M prompts from the LAION-COCO [48]. Each feature map has dimensions of 16×16 , resulting in a training dataset of 384M dense feature vectors per transformer block. For the SAE training process, we followed the methodology described in [19], using the TopK activation function and an auxiliary loss to handle dead features. For more details on the SAE training and for training metrics consider the App. E.

4. Qualitative Analysis

In this section, we perform a visual qualitative analysis to gain deeper insight into the behavior and characteristics of the learned features across transformer blocks. First, we introduce feature visualization techniques and then use them to conduct a case study. For the sake of simplicity in the notation, we omit the transformer block index ℓ .

4.1. Feature Visualization Techniques

We introduce our methods used for feature visualization used in Fig. 3. Informally, given a feature, *spatial activations* (denoted by `hmap`) highlight the regions of an image where the feature activates during generation process. *Activation modulation* (A. columns) refers to the intervention process in which the feature activations are enhanced or diminished. This technique is used to demonstrate how the manipulation of a feature’s value affects the generated image. Finally, *empty-prompt interventions* (B. column) illustrate the isolated role of the feature by disabling all other features during generation conditioned on an empty prompt. In the remainder of this section, we provide formal definitions and details.

Spatial activations. We visualize a sparse feature map $S^\rho \in \mathbb{R}^{h \times w}$ containing activations of a feature ρ across the spatial locations by up-scaling it to the size of the generated images and overlaying it as a heatmap over the generated images. In the heatmap, red indicates the highest feature activation, and blue represents the lowest non-zero one.

Top dataset examples. For a given feature ρ , we sort dataset examples according to their average spatial activation

$$a_\rho = \frac{1}{wh} \sum_{i=1}^h \sum_{j=1}^w S_{ij}^\rho \in \mathbb{R}. \quad (8)$$

We use equation 8 to define the top dataset examples and to sample from the top 5% quantile of the activating examples ($a_\rho > 0$). We will refer to them as top 5% images for a feature ρ .

Note that S_{ij}^ρ always depends on an embedding of the input prompt c and input noise z_1 , via $S_{ij}(c, z_1) = \text{ENC}(\Delta D_{ij}(c, z_1))$, which we usually omit for ease of notation. As a result, a_ρ also depends on c and z_1 . When we refer to the top dataset examples, we mean our (c, z_1) pairs with the largest values for $a_\rho(c, z_1)$.

Activation modulation. We design interventions that allow us to modulate the strength of the ρ th feature. Specifically, we achieve this by adding or subtracting a multiple of the feature ρ on all of the spatial locations i, j proportional to its original activation S_{ij}^ρ

$$\Delta D'_{ij} = \Delta D_{ij} + \beta S_{ij}^\rho \mathbf{f}_\rho, \quad (9)$$

in which ΔD_{ij} is the update performed by the transformer block before and $\Delta D'_{ij}$ after the intervention, $\beta \in \mathbb{R}$ is a

modulation factor, and \mathbf{f}_ρ is the ρ th learned feature vector. In the following, we will refer to this intervention as *activation modulation intervention*.

Note that S_{ij}^ρ can be also freely defined allowing for the application of sparse features to arbitrary images and spatial positions (refer to Fig. 1 for examples).

Activation on empty context. Another way of visualizing the causal effect of features is to activate them while doing a forward pass on the empty prompt $c(\text{“”})$. To do so, we turn off all other features at the transformer block ℓ of intervention and turn on the target feature ρ . Formally, we modify the forward pass by setting

$$D_{ij}^{\text{out}'} = D_{ij}^{\text{in}} + \gamma k \mu_\rho \mathbf{f}_\rho, \quad (10)$$

in which $D_{ij}^{\text{out}'}$ replaces residual stream plus transformer block update, D_{ij}^{in} is the input to the block, \mathbf{f}_ρ is the ρ th learned feature vector, $\gamma \in \mathbb{R}$ is a hyperparameter to adjust the intervention strength, and μ_ρ is a feature-dependent multiplier obtained by taking the average activation across positive activations of ρ (collected over a subset of 50,000 dataset examples). Multiplying it by k aims to recover the coefficients lost by setting the other features to zero. Further in the text, we will refer to this intervention as *empty-prompt intervention*, and the images generated using this method with γ set to 1, as *empty-prompt intervention images*.

Note that we directly added/subtracted feature vectors to the dense vectors for both intervention types instead of encoding, manipulating sparse features, and decoding. This approach helps mitigate side effects caused due to reconstruction loss (see App. F).

4.2. Top Features Activated by a Prompt

Combining the feature visualization techniques in Fig. 3, we depict the features with the highest average activation when processing the prompt: “A cinematic shot of a professor sloth wearing a tuxedo at a BBQ party”. We present an analysis of the transformer blocks in order of decreasing interpretability.

Down.2.1. The `down.2.1` transformer block appears to contribute to the image composition. Several features relate to the prompt: 4539 “sloth”, 4751 “a tuxedo”, 2881 “party”.

Activation modulation interventions with negative β (A. -6.0 columns) result in removing or changing scene objects in ways that align with the heatmap (`hmap` column) and the top examples (C columns): 1674 *removes* the light chains in the back, 4608 the umbrellas/tents, 4539 the 3D animation-like sloth face, and, 4751 *changes* the type of suit. Similarly, enhancing the same features (A. 6.0 column) makes the corresponding elements more visible and distinct.

Notably, activating the features on the empty prompt often creates meaningful images with related elements (B.

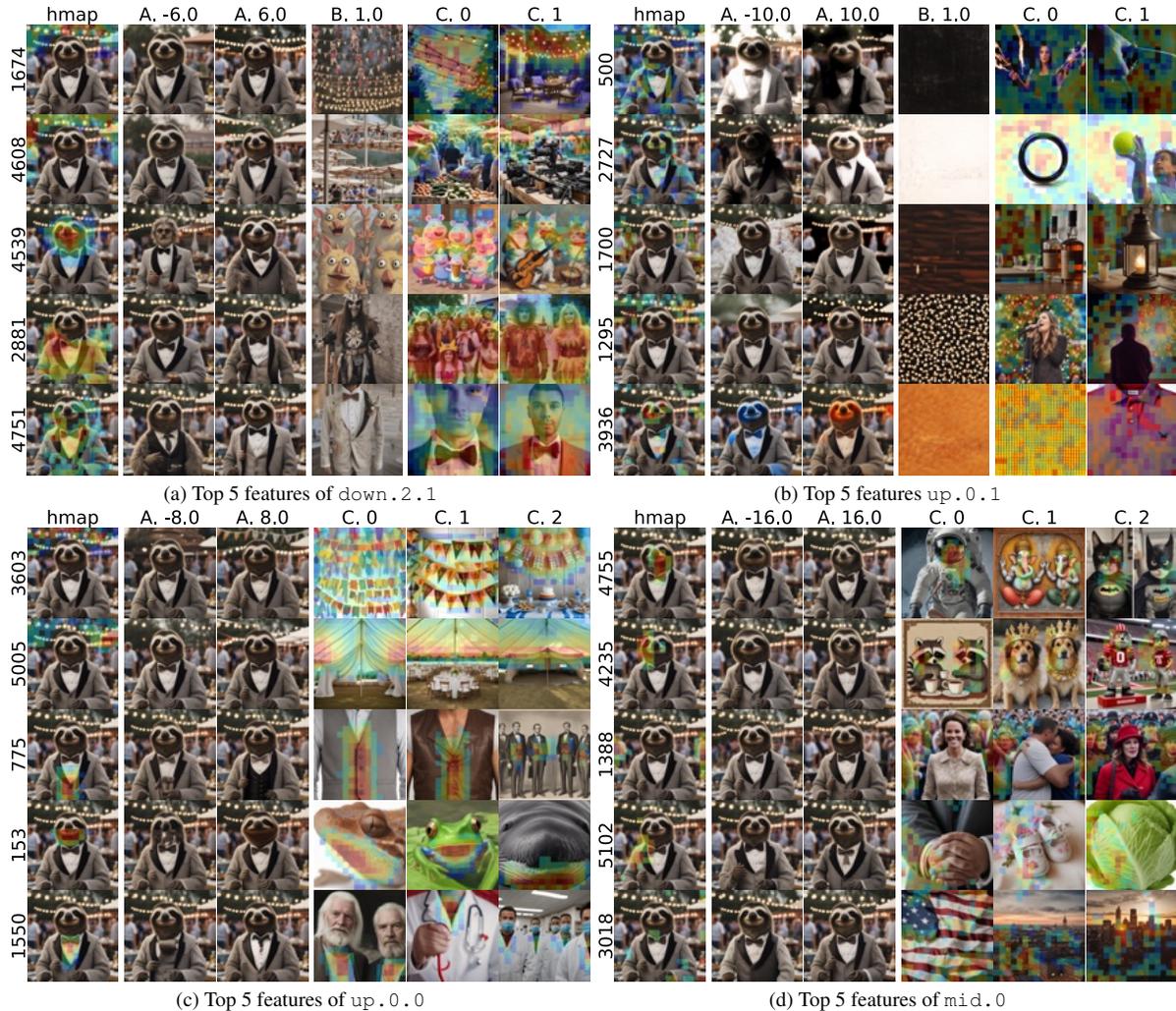


Figure 3. The top 5 features of `down.2.1` (a), `up.0.1` (b), `up.0.0` (c) and `mid.0` (d) for the prompt: "A cinematic shot of a professor sloth wearing a tuxedo at a BBQ party." Each row represents a feature. The first column depicts a feature heatmap (highest activation red and lowest nonzero one blue). The columns with titles containing "A" show feature modulation interventions, those containing "B" the intervention of turning on the feature on the empty prompt, and the ones containing "C" depict top dataset examples. Floating point values in the title denote β and γ values.

column). For reference, with the fixed random seed we use, the empty prompt generation without interventions resembles a painting of a piece of nature with a lot of green and brown tones.

While top dataset examples (C.0, C.1) and empty prompt intervention (B.) mostly agree with the feature activation heatmaps (hmap column), some of them provide additional insights, e.g., 2881, which activates on the suit, seems to correspond to (masqueraded) characters in a (festive) scene.

Up.0.1. Based on our observations, the features of `up.0.1` appear to contribute to the style.

Interestingly, turning on the `up.0.1` features on the empty prompt (B. column) results in texture-like images. Furthermore, the contribution of activation modulation in-

terventions (A. columns) is highly localized and most of inactive image area remains unchanged. For the `up.0.1` we find it remarkable that often the features' ablations and amplifications are counterparts: 500 (light, shadow), 2727 (shadow, light), 3936 (blue, orange).

Up.0.0. For the third transformer block, `up.0.0`, we observe that most top dataset examples and their activations (C columns) are quite interpretable: 3603 corresponds to party decorations, 5005 to the upper part of a tent, 775 to buttons on a suit, 153 to the lower animal jaw, 1550 to collars. All the features exhibit an expected causal effect on the generation when ablated or enhanced (A. columns).

The activation regions of the features often are very concentrated. Similarly to `up.0.1`, activation modulation in-

terventions leave inactive image regions mostly unaffected. For the empty prompt, activating these features produces abstract-looking images that are hard to relate to the other columns. Thus, we excluded this visualization technique and instead added one more dataset example. In summary, the learned features of this transformer block primarily add local details to the generation and, importantly, they are effective only within a suitable context.

Mid.0. The specific role of the fourth block (`mid.0`) is not well understood. We find it more difficult to interpret because most interventions in the `mid.0` block have subtle effects. We did not include empty-prompt intervention results because they barely affect the generation.

Despite these subtle effects, dataset examples (C. columns) and heatmaps (hmap column) mostly agree with each other and are specific enough to be interpretable: 4755 activates on bottom right part of faces, 4235 on left part of (animal) faces, 1388 on people in the background, and, 5102 on outlines the left border of the main object in the scene. We hypothesize that `mid.0`'s features are more abstract, potentially encoding spatial location and relations between objects.

Random feature visualizations. In the App. C, we provide additional visualizations of a sample of features by presenting their top activating examples and the effects of interventions.

5. Quantitative Evaluation

In this section, we follow up on qualitative insights by collecting quantitative evidence.

5.1. Annotation Pipeline

Feature annotation with an LLM followed by further evaluation is a common way to assess feature properties such as specificity, sensitivity, and causality [7]. We found it applicable to the features learned by the `down.2.1` transformer block, which have a strong effect on the generation. Thus, they are amendable to automatic annotation using visual language models (VLMs) such as GPT-4o [35]. In contrast, for the features of other blocks with more subtle effects, we found VLM-generated captions to be unsatisfactory. In order to caption the features of `down.2.1`, we prompt GPT-4o with a sequence of 14 images. The first five images are irrelevant to the feature (i.e., the feature was inactive during the generation of the images), followed by a progression of 4 images with increasing average activation values, and finished by five images with the highest average activation values. The last nine images are provided alongside their so-called ‘‘coldmaps’’: a version of an image with weakly active and inactive regions being faded and concealed. The prompt template and examples of the captions can be found in the App. G.

5.2. Experimental Details

We perform a series of experiments to get statistical insights into the features. We report the majority of the experimental scores in the format $M(S)$. When the score is reported in the context of a SDXL Turbo transformer block, it means that we computed the score for each feature of the block and set M and S to mean and standard deviation across the feature scores. Note that S does not represent the error margin of M , as the actual error margin is much lower.⁵ Therefore, almost all the differences in the reported means are statistically significant. For the baselines, we calculate the mean and standard deviation across the scores of a 100-element sample.

Table 1. Specificity, texture score, and color activation for different blocks and baselines.

Block	Specificity	Texture	Color
<code>Down.2.1</code>	0.71 (0.11)	0.16 (0.02)	86.2 (14.9)
<code>Mid</code>	0.62 (0.11)	0.14 (0.01)	84.7 (16.3)
<code>Up.0.0</code>	0.66 (0.12)	0.18 (0.03)	86.3 (16.5)
<code>Up.0.1</code>	0.65 (0.11)	0.20 (0.02)	73.8 (20.6)
Random	0.50 (0.10)	0.13 (0.02)	90.7 (54.9)
Same Prompt	0.89 (0.06)	–	–
Textures	–	0.18 (0.02)	–

Interpretability. Features are usually considered interpretable if they are sufficiently specific, i.e., images exhibiting the feature share some commonality. In order to measure this property, we compute the similarity between images on which the feature is active. High similarity between these images is a proxy for high specificity. For each feature, we collect 10 random images among top 5% images for this feature and calculate their average pairwise CLIP similarity [10, 39]. This value reflects how semantically similar the contexts are in which the feature is most active. We display the results in the first column of Table 1, which shows that the CLIP similarity between images with the feature active is significantly higher than the random baseline (CLIP similarity between random images) for all transformer blocks. This suggests that the generated images share similarities when a feature is active.

For `down.2.1` we compute an additional *interpretability* score by comparing how well the generated annotations align with the top 5% images. The resulting CLIP similarity score is 0.21 (0.03) and significantly higher than the random baseline (average CLIP similarity with random images) 0.12 (0.02). To obtain an upper bound on this score we also compute the CLIP similarity to an image generated from the feature annotation, which is 0.25 (0.03).

Causality. We can use the feature annotations to measure a feature’s causal strength by comparing the empty

⁵Given that M is computed over a sample of 1280 elements, the confidence interval of M can be estimated as $M \pm S \cdot 0.055$.

prompt intervention images with the caption.⁶ The CLIP similarity between intervention images and feature caption is 0.19 (0.04) and almost matches the annotation-based interpretability score of 0.21 (0.03). This suggests that feature annotations effectively describe to the corresponding empty-prompt intervention images. Notably, the annotation pipeline did not use empty-prompt intervention images to generate captions. This fact speaks for the high causal strength of the features learned on `down.2.1`.

Sensitivity. A feature is considered sensitive when activated in its relevant context. As a proxy for the context, we have chosen the feature annotations obtained with the auto-annotation pipeline. For each learned feature, we collected the 100 prompts from a 1.5M sample of LAION-COCO with the highest sentence similarity based on sentence transformer embeddings of `all-MiniLM-L6-v2` [43]. Next, we run SDXL Turbo on these prompts and count the proportion of generated images in which the feature is active on more than 0%, 10%, 30% of the image area, resulting in 0.60 (0.32), 0.40 (0.34), 0.27 (0.30) respectively, which is much higher than the random baseline, which is at 0.06 (0.09), 0.003 (0.006), 0.001 (0.003). However, the average scores are < 1 and thus not perfect. This may be caused by incorrect or imprecise annotations for subtle features and, therefore, hard to annotate with a VLM and SDXL Turbo failing to comply with some prompts.

Relatedness to texture. In Fig. 3 the empty prompt interventions of the `up.0.1` features resulted in texture-like pictures. To quantify whether this consistently happens, we design a simple texture score by computing CLIP similarity between an image and the word “texture”. Using this score, we compare empty-prompt interventions of the different transformer blocks with each other and real-world texture images. The results are in the second column of Table 1 and suggest that empty-prompt intervention images of `up.0.1` and `up.0.0` resemble textures and some of the `down.2.1` images look like textures as well. For `up.0.0`, we did not observe any connection of these images to the top activating images. Interestingly, the score of `up.0.1` is higher than the one of the real-world textures dataset (Cimpoi et al. [11]).

Color sensitivity. In our qualitative analysis, we suggested that the features learned on `up.0.1` relate to texture and color. If this holds, the image regions that activate a feature should not differ significantly in color on average. To test that, we calculate the “average” color for each feature: this is a weighted average of pixel colors with the feature activation values as weights. To determine the average color of a each feature we compute it over a sample of 10 images of the feature’s top 5% images. Then, we calculate Manhattan distances between the colors of the pixels and the “aver-

age” color on the same images (the highest possible distance is $3 \cdot 255 = 765$). Finally, we take a weighted average of the Manhattan distances using the same weights. We report these distances for different transformer blocks and for the images generated on random prompts from LAION-COCO. We present the results in the third column of Table 1. The average distance for the `up.0.1` transformer block is, in fact, the lowest.

Table 2. Manhattan distances between original and intervened images at varying intervention strengths outside/inside of the feature’s activation map.

Block	-10	-5	5	10
<code>Down.2.1</code>	148.2 / 116.0	124.2 / 94.4	101.4 / 78.7	128.9 / 105.60
<code>Mid</code>	69.2 / 32.2	39.4 / 18.5	33.2 / 15.2	59.9 / 29.82
<code>Up.0.0</code>	105.3 / 38.4	77.7 / 23.7	63.6 / 23.3	88.6 / 37.08
<code>Up.0.1</code>	125.0 / 26.8	73.1 / 16.4	68.6 / 21.9	98.9 / 34.74

Intervention locality. We suggested that features learned on `up.0.0` and `up.0.1` primarily influence local regions of the generation, with minimal effect outside the active areas. To test this, we measure changes in the top 5% images inside and outside the active regions while performing activation modulation interventions. To exclude weak activation regions from consideration, a pixel is considered inside the active area if the corresponding patch has an activation value larger than 50% of the image patches, and it is outside the active area if the corresponding patch has zero activation. Table 2 reports Manhattan distances between the original images and the intervened images outside and inside the active areas for activation modulation intervention strengths -10, -5, 5, 10. The features for `up.0.0` and `up.0.1` have a stronger effect inside the active area than outside, unlike `down.2.1` where the difference is smaller.

6. Related Work

In this section, we discuss relevant prior research.

Layer specialization in diffusion models. Similar to our findings on the roles of SDXL Turbo’s transformer blocks, [1, 52, 56] observe specializations among the layers and denoising steps of text-to-image diffusion models as well. Voynov et al. [52] introduce layer-specific embeddings for the text conditioning and find that different sets of layers are more effective for influencing the generation of specific concepts such as styles or objects. For a selected set of image attributes (e.g., color, object, layout, style), Agarwal et al. [1] analyze which attributes are captured in which timestamp and layer. They also find that a subset of these attributes is often captured in the same layer and across the same denoising step. Zhang et al. [56] observe that during the denoising process of diffusion models first the layout forms (in early timestamps), then the content, and finally the material and style.

⁶We require feature captions for the causality and sensitivity analyses, we only have them for `down.2.1`.

In contrast to these prior works, our approach takes a fundamentally different perspective and does not rely on either handcrafted attributes or prompts. Instead, we introduce a new lens for analyzing SDXL Turbo’s transformer blocks, which reveals specialization among the blocks as well. Interestingly, our findings on SDXL Turbo, which is a distilled, few-step diffusion model parallel [56]’s observations, which identify that composition precede material and style. In SDXL Turbo, this progression occurs across the layers instead of the denoising timestamps.

Analyzing the latent space of diffusion models. Kwon et al. [27] show that diffusion models have a semantically meaningful latent space. Park et al. [36] analyze the latent space of diffusion models using Riemannian geometry. Li et al. [28] and Dalva and Yanardag [13] present self-supervised methods for finding semantic directions. Similarly, Gandikota et al. [18] show that the attribute variations lie in a low-rank space by learning LoRA adapters [22] on top of pre-trained diffusion models. Brack et al. [4] and Wang et al. [53] demonstrate effective semantic vector algebraic operations in the latent space of DMs, as observed by Mikolov et al. [33]. However, none of those works train SAEs to interpret and control the latent space.

Mechanistic interpretability. Sparse autoencoders have recently been popularized by [5], in which they show that it is possible to learn interpretable features by decomposing neuron activations in MLPs in 2-layer transformer language models. At the same time, a parallel work decomposed the elements of the residual stream [12], which followed up on [49]. To our knowledge, the first work that applied sparse autoencoders to transformer-based LLM was [54], which learned a joint dictionary for features of all layers. Recently, sparse autoencoders have gained much traction, and many have been trained even on state-of-the-art LLMs [19, 29, 51]. In addition, great tools are available for inspection [30] and automatic interpretation [7] of learned features. [32] have shown how to use SAE features to facilitate automatic circuit discovery.

The studies most closely related to our work are [2], [23] and [14]. Ismail et al. [23] apply concept bottleneck methods [25] that decompose latent concepts into vectors of interpretable concepts to generative image models, including diffusion models. Unlike the SAEs that we train, this method requires labeled concept data. Daujotas [14] decomposes CLIP [10, 39] vision embeddings using SAEs and use them for conditional image generation with a diffusion model called Kandinsky [42]. Importantly, using SAE features, they are able to manipulate the image generation process in interpretable ways. In contrast, in our work, we train SAEs on intermediate representations of the forward pass of SDXL Turbo. Consequently, we can interpret and manipulate SDXL Turbo’s forward pass on a finer granularity, e.g., by intervening on specific transformer blocks and

spatial positions. Another closely related work to ours is [2], in which neurons in generative adversarial neural networks are interpreted and manipulated. The interventions in [2] are similar to ours, but on neurons instead of sparse features. In order to identify neurons corresponding to a semantic concept, [2] require semantic segmentation maps.

7. Conclusion and Discussion

We trained SAEs on SDXL Turbo’s opaque intermediate representations. This study is the first in the academic literature to mechanistically interpret the intermediate representations of a modern text-to-image model. Our findings demonstrate that SAEs can extract interpretable features and have a significant causal effect on the generated images. Importantly, the learned features provide insights into SDXL Turbo’s forward pass, revealing that transformer blocks fulfill specific and varying roles in the generation process. In particular, our results clarify the functions of `down.2.1`, `up.0.0`, and `up.0.1`. However, the role of `mid.0` remains less defined; it seems to encode more abstract information and interventions are less effective.

We follow up with a discussion of the results and their implications for future research. Based on our observations, we suggest a preliminary hypothesis about SDXL Turbo’s generation process: `down.2.1` decides on top-level composition, `mid.0` encodes low-level semantics, `up.0.0` adds details based on the two above, and `up.0.1` fills in color, texture, and style.

Our work focuses on analyzing SDXL Turbo’s intermediate representations. As a relatively compact, few-step diffusion model with a small number of naturally partitioned components, SDXL Turbo turned out to be convenient to analyze with SAEs. However, the application of the proposed techniques to larger and more complex text-to-image diffusion models with alternative architectures represents a promising direction for further research. We provide some motivational results on PixArt LCM [9] in the App. H. Additionally, we observe that SAE features learned on SDXL Turbo’s one-step generation are applicable to modify few-step (2-4) generation workflow (App. D).

Analyzing larger diffusion models with higher number of diffusion steps would benefit from advanced interpretability techniques capable of capturing connections between its components and across denoising steps. For example, such techniques are explored in [31, 32]. Marks et al. [32] compute circuits showing how different layers and attention heads wire together and Lindsey et al. [31] introduce cross-coders, a variation of SAEs that allows to learn a shared set of features over latents corresponding to different layers.

Our work highlights the potential of SAEs in revealing the internal structure of diffusion models like SDXL Turbo, and it could help future researchers answer more sophisticated questions about image generation.

Acknowledgements

The authors thank Danila Zubko for the initial contribution, David Bau, Tim R. Davidson, Niv Cohen, Gytis Daujotas and Alexander Sharipov for the valuable discussions and feedback, as well as RCP and IC clusters maintainers.

Robert West’s lab is partly supported by grants from the Swiss National Science Foundation (200021_185043, TMSGI2_211379), Swiss Data Science Center (P22_08), H2020 (952215), Microsoft, and Google.

Çaglar Gulcehre’s lab is supported by nimble.ai.

References

- [1] Aishwarya Agarwal, Srikrishna Karanam, Tripti Shukla, and Balaji Vasanth Srinivasan. An image is worth multiple words: Multi-attribute inversion for constrained text-to-image synthesis. *arXiv preprint arXiv:2311.11919*, 2023. 7
- [2] David Bau, Jun-Yan Zhu, Hendrik Strobelt, Bolei Zhou, Joshua B Tenenbaum, William T Freeman, and Antonio Torralba. Gan dissection: Visualizing and understanding generative adversarial networks. In *International Conference on Learning Representations*, 2019. 8
- [3] James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jianfeng Wang, Linjie Li, Long Ouyang, Juntang Zhuang, Joyce Lee, Yufei Guo, et al. Improving image generation with better captions. *Computer Science*. <https://cdn.openai.com/papers/dall-e-3.pdf>, 2(3):8, 2023. 1
- [4] Manuel Brack, Felix Friedrich, Dominik Hintersdorf, Lukas Struppek, Patrick Schramowski, and Kristian Kersting. Sega: Instructing text-to-image models using semantic guidance. In *Advances in Neural Information Processing Systems*, pages 25365–25389. Curran Associates, Inc., 2023. 8
- [5] Trenton Bricken, Adly Templeton, Joshua Batson, Brian Chen, and Adam Jermy et al. Towards monosemanticity: Decomposing language models dictionary learning. *Transformer Circuits*, 2023. 1, 2, 8, 12, 13
- [6] Tom B Brown. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020. 1
- [7] Juang Caden, Paulo Gonalo, Drori Jacob, and Belrose Nora. Open source automated interpretability for sparse autoencoder features, 2024. Accessed: 2024-09-27. 2, 6, 8
- [8] Haozhe Chen, Carl Vondrick, and Chengzhi Mao. Selfie: Self-interpretation of large language model embeddings. *arXiv preprint arXiv:2403.10949*, 2024. 1
- [9] Junsong Chen, Yue Wu, Simian Luo, Enze Xie, Sayak Paul, Ping Luo, Hang Zhao, and Zhenguo Li. Pixart- $\{\delta\}$: Fast and controllable image generation with latent consistency models. *arXiv preprint arXiv:2401.05252*, 2024. 8, 15
- [10] Mehdi Cherti, Romain Beaumont, Ross Wightman, Mitchell Wortsman, Gabriel Ilharco, Cade Gordon, Christoph Schuhmann, Ludwig Schmidt, and Jenia Jitsev. Reproducible scaling laws for contrastive language-image learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2818–2829, 2023. 6, 8
- [11] M. Cimpoi, S. Maji, I. Kokkinos, S. Mohamed, and A. Vedaldi. Describing textures in the wild. In *Proceedings of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2014. 7
- [12] Hoagy Cunningham, Aidan Ewart, Logan Riggs, Robert Huben, and Lee Sharkey. Sparse autoencoders find highly interpretable features in language models. *arXiv preprint arXiv:2309.08600*, 2023. 1, 2, 8, 12
- [13] Yusuf Dalva and Pinar Yanardag. Noiseclr: A contrastive learning approach for unsupervised discovery of interpretable directions in diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 24209–24218, 2024. 8
- [14] Gytis Daujotas. Interpreting and steering features in images, 2024. Accessed: 2024-09-27. 8
- [15] Jacob Devlin. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018. 1
- [16] Prafulla Dhariwal and Alex Nichol. Diffusion models beat gans on image synthesis, 2021. 1
- [17] Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec, Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, Roger Grosse, Sam McCandlish, Jared Kaplan, Dario Amodei, Martin Wattenberg, and Christopher Olah. Toy models of superposition. *Transformer Circuits*, 2022. 2, 13
- [18] Rohit Gandikota, Joanna Materzynska, Tingrui Zhou, Antonio Torralba, and David Bau. Concept sliders: Lora adaptors for precise control in diffusion models, 2023. 8
- [19] Leo Gao, Tom Dupr e la Tour, Henk Tillman, Gabriel Goh, Rajan Troll, Alec Radford, Ilya Sutskever, Jan Leike, and Jeffrey Wu. Scaling and evaluating sparse autoencoders. *arXiv preprint arXiv:2406.04093*, 2024. 3, 8, 13
- [20] Asma Ghandeharioun, Avi Caciularu, Adam Pearce, Lucas Dixon, and Mor Geva. Patchscope: A unifying framework for inspecting hidden representations of language models. *arXiv preprint arXiv:2401.06102*, 2024. 1
- [21] Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022. 1
- [22] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models, 2021. 8
- [23] Aya Abdelsalam Ismail, Julius Adebayo, Hector Corrada Bravo, Stephen Ra, and Kyunghyun Cho. Concept bottleneck generative models. In *The Twelfth International Conference on Learning Representations*, 2023. 8
- [24] William B Johnson, Joram Lindenstrauss, and Gideon Schechtman. Extensions of lipschitz maps into banach spaces. *Israel Journal of Mathematics*, 54(2):129–138, 1986. 13
- [25] Pang Wei Koh, Thao Nguyen, Yew Siang Tang, Stephen Mussmann, Emma Pierson, Been Kim, and Percy Liang. Concept bottleneck models. In *International conference on machine learning*, pages 5338–5348. PMLR, 2020. 8

- [26] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25, 2012. 1
- [27] Mingi Kwon, Jaeseok Jeong, and Youngjung Uh. Diffusion models already have a semantic latent space, 2023. 8
- [28] Hang Li, Chengzhi Shen, Philip Torr, Volker Tresp, and Jindong Gu. Self-discovering interpretable diffusion latent directions for responsible text-to-image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12006–12016, 2024. 8
- [29] Tom Lieberum, Senthoooran Rajamanoharan, Arthur Conmy, Lewis Smith, Nicolas Sonnerat, Vikrant Varma, János Kramár, Anca Dragan, Rohin Shah, and Neel Nanda. Gemma scope: Open sparse autoencoders everywhere all at once on gemma 2. *arXiv preprint arXiv:2408.05147*, 2024. 8
- [30] Johnny Lin and Joseph Bloom. Neuronpedia: Interactive reference and tooling for analyzing neural networks, 2023. Software available from neuronpedia.org. 8
- [31] Jack Lindsey, Adly Templeton, Jonathan Marcus, Thomas Conerly, Joshua Batson, and Christopher Olah. Sparse cross-coders for cross-layer features and model diffing. *Transformer Circuits*, 2024. 8
- [32] Samuel Marks, Can Rager, Eric J Michaud, Yonatan Belinkov, David Bau, and Aaron Mueller. Sparse feature circuits: Discovering and editing interpretable causal graphs in language models. *arXiv preprint arXiv:2403.19647*, 2024. 2, 8
- [33] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space, 2013. 8
- [34] Chris Olah, Nick Cammarata, Ludwig Schubert, Gabriel Goh, Michael Petrov, and Shan Carter. An overview of early vision in inceptionv1. *Distill*, 5(4):e00024–002, 2020. 1
- [35] OpenAI. Hello gpt-4o, 2024. Accessed: 2024-09-28. 6
- [36] Yong-Hyun Park, Mingi Kwon, Jaewoong Choi, Junghyo Jo, and Youngjung Uh. Understanding the latent space of diffusion models through the lens of riemannian geometry. In *Advances in Neural Information Processing Systems*, pages 24129–24142. Curran Associates, Inc., 2023. 8
- [37] Pablo Pernias, Dominic Rampas, Mats L Richter, Christopher J Pal, and Marc Aubreville. Würstchen: An efficient architecture for large-scale text-to-image diffusion models. *arXiv preprint arXiv:2306.00637*, 2023. 1
- [38] 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, 2023. 1, 3
- [39] 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 *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 6, 8
- [40] Daking Rai, Yilun Zhou, Shi Feng, Abulhair Saparov, and Ziyu Yao. A practical review of mechanistic interpretability for transformer-based language models. *arXiv preprint arXiv:2407.02646*, 2024. 1
- [41] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International conference on machine learning*, pages 8821–8831. Pmlr, 2021. 1
- [42] Anton Razzhigaev, Arseniy Shakhmatov, Anastasia Maltseva, Vladimir Arkhipkin, Igor Pavlov, Ilya Ryabov, Angelina Kuts, Alexander Panchenko, Andrey Kuznetsov, and Denis Dimitrov. Kandinsky: An improved text-to-image synthesis with image prior and latent diffusion. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 286–295, 2023. 8
- [43] Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 2019. 7
- [44] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models, 2022. 1, 3
- [45] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in neural information processing systems*, 35:36479–36494, 2022.
- [46] Axel Sauer, Dominik Lorenz, Andreas Blattmann, and Robin Rombach. Adversarial diffusion distillation, 2023. 1, 3
- [47] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *Advances in Neural Information Processing Systems*, 35:25278–25294, 2022. 1, 2
- [48] Christoph Schuhmann, Andreas Köpf, Richard Vencu, Theo Coombes, Romain Beaumont, and Benjamin Trom. Laion coco: 600m synthetic captions from laion2b-en, 2022. Accessed: 2024-10-01. 2, 3
- [49] Lee Sharkey, Dan Braun, and beren. Interim research report: Taking features out of superposition with sparse autoencoders, 2022. Accessed: 2024-09-27. 8
- [50] Matteo Spinelli. Advanced style transfer with the mad scientist node. YouTube video, 2024. Accessed: 2024-09-17. 2, 11
- [51] Adly Templeton and Tom Conerly et al. Scaling monosemanticity: Extracting interpretable features from claude 3 sonnet, 2024. Accessed: 2024-09-27. 8
- [52] Andrey Voynov, Qinghao Chu, Daniel Cohen-Or, and Kfir Aberman. p+: Extended textual conditioning in text-to-image generation. *arXiv preprint arXiv:2303.09522*, 2023. 7
- [53] Zihao Wang, Lin Gui, Jeffrey Negrea, and Victor Veitch. Concept algebra for (score-based) text-controlled generative

models. In *Advances in Neural Information Processing Systems*, pages 35331–35349. Curran Associates, Inc., 2023. 8

- [54] Zeyu Yun, Yubei Chen, Bruno A Olshausen, and Yann LeCun. Transformer visualization via dictionary learning: contextualized embedding as a linear superposition of transformer factors. *arXiv preprint arXiv:2103.15949*, 2021. 2, 8, 12
- [55] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018. 14
- [56] Yuxin Zhang, Weiming Dong, Fan Tang, Nisha Huang, Haibin Huang, Chongyang Ma, Tong-Yee Lee, Oliver Deussen, and Changsheng Xu. Prospect: Prompt spectrum for attribute-aware personalization of diffusion models. *ACM Transactions on Graphics (TOG)*, 42(6):1–14, 2023. 7, 8

A. Finding Causally Influential Transformer Blocks

We narrow down design space of the 11 cross-attention transformer blocks (see Fig. 4) to those with the highest causal impact on the output. In order to assess their causal impact on the output we qualitatively study the effect of individually ablating each of them (see Fig. 5). As can be seen in Fig. 5 each of the middle blocks `down.2.1`, `mid.0`, `up.0.0`, `up.0.1` have a relatively high impact on the output respectively. In particular, the blocks `down.2.1` and `up.0.1` stand out. It seems like most colors and textures are added in `up.0.1`, which in the community is already known as “style” block [50]. Ablating `down.2.1`, which is also already known in the community as “composition” block, impacts the entire image composition, including object sizes, orientations and framing. The effects of ablating other blocks such as `mid.0` and `up.0.0` are more subtle. For `mid.0` it is difficult to describe in words and `up.0.0` seems to add local details to the image while leaving the overall composition mostly intact.

B. Case Study: Most Active Features on a Prompt

Fig. 10 demonstrates an extended version of qualitative analysis case study, showcasing 9 top features instead of 5.

C. Case Study: Random Features

In this case study, we explore the learned features independently of any specific prompt. In Fig. 11 and Fig. 12, we demonstrate the first 5 and last 5 learned features for each transformer block (since SAEs were initialized randomly before training, we can treat these features as a random sample). As SAEs are randomly initialized before the training process, these sets can be considered as random samples of features. Each feature visualization consists of 3 images

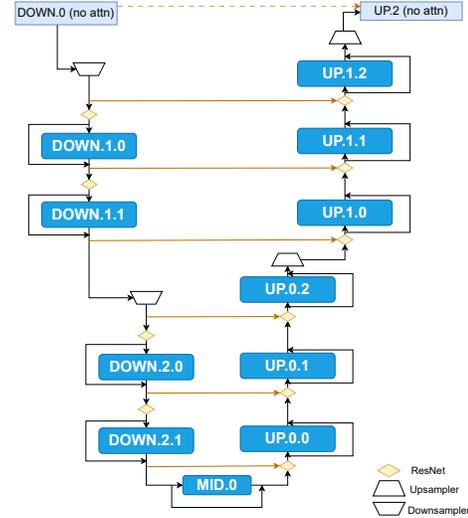


Figure 4. Cross-attention transformer blocks in SDLX’s U-net.

of top 5% images for this feature, and their perturbations with activation modulation interventions. For `down.2.1` and `up.0.1`, we also include the empty-prompt intervention images. Additionally, we provide visualizations of several selected features in App. C Fig. 13 and demonstrate the effects of their forced activation on unrelated prompts in App. C Fig. 14.

Feature plots. We provide the same plots as in Fig. 11 but for the last six feature indices of each transformer block in Fig. 12 and the corresponding prompts in Table 6. Additionally, provide some selected features for `down.2.1` and `up.0.1` in Fig. 13 and the corresponding prompts in Table 7.

Intervention plots. Additionally, we provide plots in which we turn on features from Fig. 13 but in unrelated prompts (as opposed to top dataset example prompts that already activate the features by themselves). For simplicity here we simply turn on the features across all spatial locations, which does not seem to be a well suitable strategy for `up.0.1`, which usually acts locally. To showcase, the difference we created one example image in Fig. 8, in which we manually draw localized masks to turn on the corresponding features.

D. Interventions on a Few-Step Generation Process

Figure 6 illustrates interventions similar to those highlighted in the main text head illustration, applied to the 4-step generation process. For this experiment, we employed equal activation maps and consistent strengths across all four denoising steps. Remarkably, the features’ functions

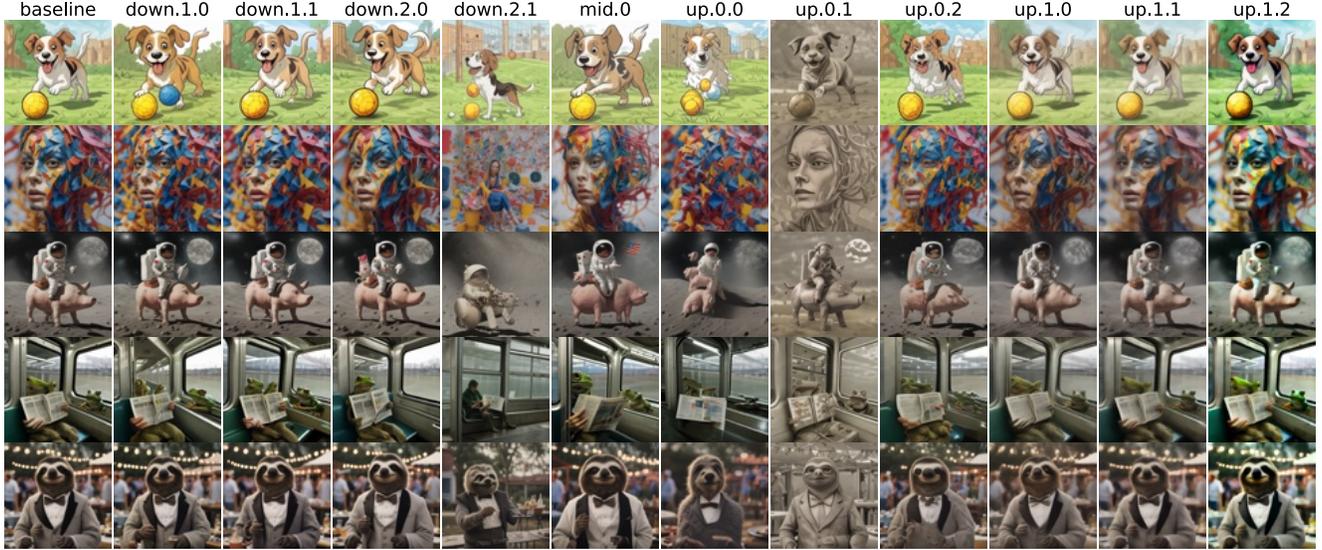


Figure 5. We generate images for the prompts “A dog playing with a ball cartoon.”, “A photo of a colorful model.”, “An astronaut riding on a pig on the moon.”, “A photograph of the inside of a subway train. There are frogs sitting on the seats. One of them is reading a newspaper. The window shows the river in the background.” and “A cinematic shot of a professor sloth wearing a tuxedo at a BBQ party.” while ablating the updates performed by different cross-attention layers (indicated by the titles). The title “baseline” corresponds to the generation without interventions.



Figure 6. Enabling features learned by sparse autoencoders in SDXL Turbo’s 4-step generation process.

were preserved, and the interventions produced images of decent quality. This indicates that features learned during the 1-step generation process may generalize effectively to a few-step setting.

E. Sparse Autoencoders and Superposition

This is an extended version of Sparse Autoencoders subsection of background section.

Let $h(x) \in \mathbb{R}^d$ be some intermediate result of a forward pass of a neural network on the input x . In a fully connected neural network, the components $h(x)$ could correspond to neurons. In transformers, which are residual neural networks with attention and fully connected layers, $h(x)$ usually either refers to the content of the residual stream after some layer, an update to the residual stream by some layer, or the neurons within a fully connected block. In general, $h(x)$ could refer to anything, e.g., also keys, queries, and values. It has been shown [5, 12, 54] that in many neural networks, especially LLMs, intermediate represen-

tations can be well approximated by sparse sums of $n_f \in \mathbb{N}$ learned feature vectors, i.e.,

$$h(x) \approx \sum_{\rho=1}^{n_f} s_{\rho}(x) \mathbf{f}_{\rho}, \quad (11)$$

where $s_{\rho}(x)$ are the input-dependent⁷ coefficients most of which are equal to zero and $\mathbf{f}_1, \dots, \mathbf{f}_{n_f} \in \mathbb{R}^d$ is a learned dictionary of feature vectors.

Importantly, these learned features are usually highly *interpretable* (specific), *sensitive* (fire on the relevant contexts), *causal* (change the output in expected ways in intervention) and usually do not correspond directly to individual neurons. There are also some preliminary results on the universality of these learned features, i.e., that different training runs on similar data result in the corresponding models picking up largely the same features [5].

Superposition. By associating task-relevant features with directions in \mathbb{R}^d instead of individual components of

⁷In the literature this input dependence is usually omitted.

$h(x) \in \mathbb{R}^d$, it is possible to represent many more features than there are components, i.e., $n_f \gg d$. As a result, in this case, the learned dictionary vectors $\mathbf{f}_1, \dots, \mathbf{f}_{n_f}$ cannot be orthogonal to each other, which can lead to interference when too many features are on (thus the sparsity requirement). However, it would be theoretically possible to have exponentially (in d) many almost orthogonal directions embedded in \mathbb{R}^d .⁸

Using representations like this, the optimization process during training can trade off the benefits of being able to represent more features than there are components in h with the costs of features interfering with each other. Such representations are especially effective if the real features underlying the data do not co-occur with each other too much, that is, they are sparse. In other words, in order to represent a single input (“Michael Jordan”) only a small subset of the features (“person”, ..., “played basketball”) is required [5, 17].

The phenomenon of neural networks that exploit representations with more features than there are components (or neurons) is called superposition [17]. Superposition can explain the presence of polysemantic neurons. The neurons, in this case, are simply at the wrong level of abstraction. The closest feature vector can change when varying a neuron, resulting in the neuron seemingly reacting to or steering semantically unrelated things.

Sparse autoencoders. In order to implement the sparse decomposition from equation 11, the vector s containing the n_f coefficients of the sparse sum is parameterized by a single linear layer followed by an activation function, called the *encoder*,

$$s = \text{ENC}(h) = (W^{\text{ENC}}(h - b_{\text{pre}}) + b_{\text{act}}), \quad (12)$$

in which $h \in \mathbb{R}^d$ is the latent that we aim to decompose, $\sigma(\cdot)$ is an activation function, $W^{\text{ENC}} \in \mathbb{R}^{n_f \times d}$ is a learnable weight matrix and b_{pre} and b_{act} are learnable bias terms. We omitted the dependencies $h = h(x)$ and $s = s(h)$ that are clear from context.

Similarly, the learnable features are parametrized by a single linear layer, called *decoder*,

$$h' = \text{DEC}(s) = W^{\text{DEC}}s + b_{\text{pre}}, \quad (13)$$

in which $W^{\text{DEC}} = (\mathbf{f}_1 | \dots | \mathbf{f}_{n_f}) \in \mathbb{R}^{d \times n_f}$ is a learnable matrix of whose columns take the role of learnable features and b_{pre} is a learnable bias term.

Training. The pair ENC and DEC are trained in a way that ensures that h' is a sparse sum of feature vectors. Given a dataset of latents h_1, \dots, h_n , both encoder and decoder

⁸It follows from the Johnson-Lindenstrauss Lemma [24] that one can find at least $\exp(d\epsilon^2/8)$ unit vectors in \mathbb{R}^d with the dot product between any two not larger than ϵ .

are trained jointly to minimize a proxy to the loss

$$\min_{W^{\text{ENC}}, W^{\text{DEC}}, b_{\text{pre}}, b_{\text{act}}} \sum_{i=1}^n \|h'_i - h_i\|_2^2 + \lambda \|s_i\|_0, \quad (14)$$

where $h_i = h(x_i)$, $s_i = \text{ENC}(h(x_i))$ (when we refer to components of s we use s_ρ instead), the $\|h'_i - h_i\|_2^2$ is a reconstruction loss, $\|s_i\|_0$ a regularization term ensuring the sparsity of the activations and λ the corresponding trade-off term.

In practice, $\|s_i\|_0$ cannot be efficiently optimized directly, which is why it is usually replaced with $\|s_i\|_1$ or other proxy objectives.

Technical details. In our work, we make use of the top- k formulation from [19], in which $\|s_i\|_0 \leq k$ is ensured by introducing the a top- k function TOPK into the encoder:

$$s = \text{ENC}(h) = \text{RELU}(\text{TOPK}(W^{\text{ENC}}(h - b_{\text{pre}}) + b_{\text{act}})). \quad (15)$$

As the name suggests, TOPK returns a vector that sets all components except the top k ones to zero.

In addition [19] use an auxiliary loss to handle dead features. During training, a sparse feature ρ is considered *dead* if s_ρ remains zero over the last 10M training examples.

The resulting training loss is composed of two terms: the L_2 -reconstruction loss and the top-auxiliary L_2 -reconstruction loss for dead feature reconstruction. For a single latent h , the loss is defined

$$L(h, h') = \|h - h'\|_2^2 + \alpha \|h - h'_{\text{aux}}\|_2^2 \quad (16)$$

In this equation, the h'_{aux} is the reconstruction based on the top k_{aux} dead features. This auxiliary loss is introduced to mitigate the issue of dead features. After the end of the training process, we observed none of them. Following [19], we set $\alpha = \frac{1}{32}$ and $k_{\text{aux}} = 256$, performed tied initialization of encoder and decoder, normalized decoder rows after each training step. The number of learned features n_f is set to 5120, which is four times the length of the input vector. The value of k is set to 10 as a good trade-off between sparsity and reconstruction quality. Other training hyperparameters are batch size: 4096, optimizer: Adam with learning rate: 10^{-4} and betas: (0.9, 0.999).

F. SAE Training Results

We trained several SAEs with different sparsity levels and sparse layer sizes and observed no dead features. To assess reconstruction quality, we processed 100 random LAION-COCO prompts through a one-step SDXL Turbo process, replacing the additive component of the corresponding transformer block with its SAE reconstruction.

The explained variance ratio and the output effects caused by reconstruction are shown in Table 3. Fig. 7

presents random examples of reconstructions from an SAE with the following hyperparameters: $k = 10, n_f = 5120$, trained on `down.2.1`. The reconstruction causes minor deviations in the images, and the fairly low LPIPS [55] and pixel distance scores also support these findings. However, to prevent these minor reconstruction errors from affecting our analysis of interventions, we decided to directly add or subtract learned directions from dense feature maps.

Table 3. Distances and explained variance ratio in generated images. “Mean” represents the average pixel Manhattan distance between original and reconstruction-intervened images, with a maximum possible value of 765. “Median” represents the median Manhattan distance per pixel, averaged over all images. ‘LPIPS’ refers to the average LPIPS score, measuring perceptual similarity. “Explained variance ratio” denotes the ratio of variance explained by the trained SAEs to the total variance.

k	n_f	Configuration	Mean	Median	LPIPS	EV (%)
5	640	down.2.1	83.29	50.04	0.3383	56.0
		mid.0	52.64	26.82	0.2032	43.4
		up.0.0	55.89	30.69	0.2276	44.8
		up.0.1	52.67	34.53	0.2073	50.3
	5120	down.2.1	74.68	41.49	0.3036	67.8
		mid.0	48.82	24.60	0.1845	50.8
		up.0.0	49.19	25.86	0.1969	57.2
		up.0.1	47.50	31.11	0.1775	59.5
10	640	down.2.1	73.65	41.79	0.2893	62.8
		mid.0	46.80	23.10	0.1772	51.5
		up.0.0	48.43	25.80	0.1908	52.5
		up.0.1	43.06	26.85	0.1638	58.7
	5120	down.2.1	64.97	34.77	0.2582	73.7
		mid.0	44.02	21.72	0.1627	58.8
		up.0.0	42.08	21.54	0.1624	64.2
		up.0.1	39.77	24.84	0.1453	67.1
20	640	down.2.1	59.29	31.47	0.2291	69.9
		mid.0	39.95	19.44	0.1459	60.0
		up.0.0	40.15	21.06	0.1499	60.9
		up.0.1	31.97	18.15	0.1196	66.7
	5120	down.2.1	56.37	29.04	0.2190	78.8
		mid.0	37.28	17.82	0.1328	66.5
		up.0.0	35.73	18.03	0.1302	70.6
		up.0.1	30.31	17.22	0.1104	74.2

G. Annotation Pipeline Details

We used GPT-4o to caption learned features on `down.2.1`. For each feature, the model was shown a series of 5 unrelated images, a progression of 9 images, the i -th of those corresponds to $\sim i \cdot 10\%$ average activation value of the maximum. Finally, we show 5 images corresponding to the highest average activations. Since some features are active on particular parts of images, the last 9 images are provided alongside their so-called “coldmaps”: a version of an image with weakly active and inactive regions being faded and concealed.

The images were generated by 1-step SDXL Turbo diffusion process on 50’000 random prompts of LAION-COCO dataset.

G.1. Textual Prompt Template

Here is the prompt template for the VLM.

System. You are an experienced mechanistic interpretability researcher that is labeling features from the hidden representations of an image generation model.

User. You will be shown a series of images generated by a machine learning model. These images were selected because they trigger a specific feature of a sparse auto-encoder, trained to detect hidden activations within the model. This feature can be associated with a particular object, pattern, concept, or a place on an image. The process will unfold in three stages:

1. **Reference Images:** First, you’ll see several images *unrelated* to the feature. These will serve as a reference for comparison.

2. **Feature-Activating Images:** Next, you’ll view images that activate the feature with varying strengths. Each of these images will be shown alongside a version where non-activated regions are masked out, highlighting the areas linked to the feature.

3. **Strongest Activators:** Finally, you’ll be presented with the images that most strongly activate this feature, again with corresponding masked versions to emphasize the activated regions.

Your task is to carefully examine all the images and identify the thing or concept represented by the feature. Here’s how to provide your response:

- **Reasoning:** Between `<thinking>` and `</thinking>` tags, write up to 400 words explaining your reasoning. Describe the visual patterns, objects, or concepts that seem to be consistently present in the feature-activating images but not in the reference images.

- **Expression:** Afterward, between `<answer>` and `</answer>` tags, write a concise phrase (no more than 15 words) that best captures the common thing or concept across the majority of feature-activating images.

Note that not all feature-activating images may perfectly align with the concept you’re describing, but the images with stronger activations should give you the clearest clues. Also pay attention to the masked versions, as they highlight the regions most relevant to the feature.

User. These images are not related to the feature: {Reference Images}

User. This is a row of 9 images, each illustrating increasing levels of feature activation. From



Figure 7. Images generated from 10 random prompts taken from the LAION-COCO dataset are shown in the first row. In the second row, down . 2 . 1 updates are replaced by their SAE reconstructions ($k = 10, n_f = 5120$). The third row visualizes the differences between the original and reconstructed images.

left to right, each image shows a progressively higher activation, starting with the image on the far left where the feature is activated at 10% relative to the image that activates it the most, all the way to the far right, where the feature activates at 90% relative to the image that activates it the most. This gradual transition highlights the feature’s growing importance across the series. {Feature-Activating Images}

User. This row consists of 9 masked versions of the original images. Each masked image corresponds to the respective image in the activation row. Areas where the feature is not activated are completely concealed by a white mask, while regions with activation remain visible. {Feature-Activating Images Coldmaps}

User. These images activate the feature most strongly. {Strongest Activators}

User. These masked images highlight the activated regions of the images that activate the feature most strongly. The masked images correspond to the images above. The unmasked regions are the ones that activate the feature. {Strongest Activators Coldmaps}

G.2. Example of Prompt Images

The images used to annotate feature 0 are shown in Fig. 9.

G.3. Examples of Generated Captions

We present the captions generated by GPT-4o for the first and last 10 features in Table 4.

H. SAE Features on Pixart-LCM XL

We trained sparse autoencoders on updates from three transformer blocks of the Pixart-LCM XL model [9] during 4-

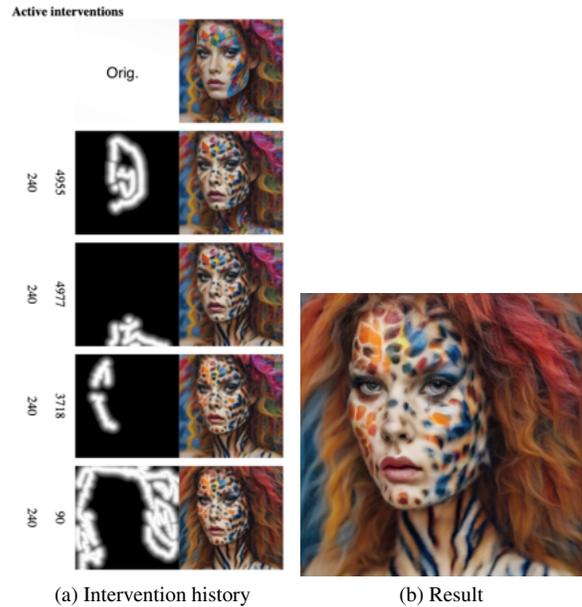


Figure 8. Local edits showcase up . 0 . 1’s ability to locally change textures in the image without affecting the remaining image. Multiple consecutive interventions are possible (a). The first in (a) row depicts the original image and each subsequent row we add an intervention by drawing a heatmap with a brush tool and then turning on the feature labelling the row only on that area. The other number (240) is the absolute feature strength of the edit. Figure (b) shows the final result in full resolution (512x512).

step generation process. Figure 15 shows activation maps of several features that were among the most active during a forward pass. We found that some features are triggered by specific objects, while others respond to particular spatial positions in the image, or their activation regions may be scattered. However, individual features typically do not exhibit the expected causal effects. We hypothesize that

Table 4. down . 2 . 1 first 10 and last 10 feature captions.

Block	Feature	Caption
down . 2 . 1	0	Organizational/storage items for documents and office supplies
	1	Luxury kitchen interiors and designs
	2	Architectural Landmarks and Monumental Buildings
	3	Upper body clothing and attire
	4	Rustic or Natural Wooden Textures or Surfaces
	5	Intricately designed and ornamental brooches
	6	Technical diagrams and instructional content
	7	Feature predominantly activated by visual representations of dresses
	8	Home decor textiles focusing on cushions and pillows
	9	Eyewear: glasses and sunglasses
	5110	Concept of containment or organized enclosure
	5111	Groups of people in collective settings
	5112	Modern minimalist interior design
	5113	Indoor plants and greenery
	5114	Feature sensitivity focused on sneakers
	5115	Handling or manipulating various objects
	5116	Athletic outerwear, particularly zippered sporty jackets
	5117	Spectator Seating in Sporting Venues
	5118	Textiles and clothing materials, focus on textures and folds
	5119	Yarn and Knitting Textiles

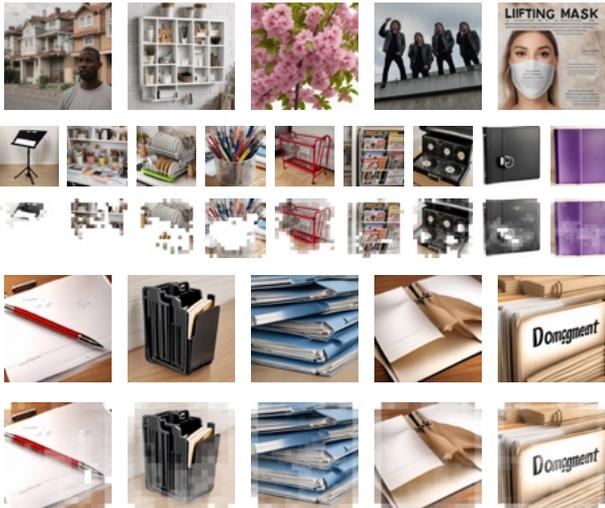


Figure 9. The images used by GPT-4o to generate captions for feature 0. From top to bottom: irrelevant images to feature 0; image progression from left to right, showing increasing activation of SAE feature 0, with low activation on the left and high activation on the right; “Coldmaps” representing the image progression; images corresponding to the highest activation of feature 0; “Coldmaps” corresponding to these highest activation images.

this may be because features across different transformer blocks are not interdependent and may exhibit causal effects only in combination. Future research is encouraged to employ advanced interpretability techniques to better understand and resolve these complex feature interactions.

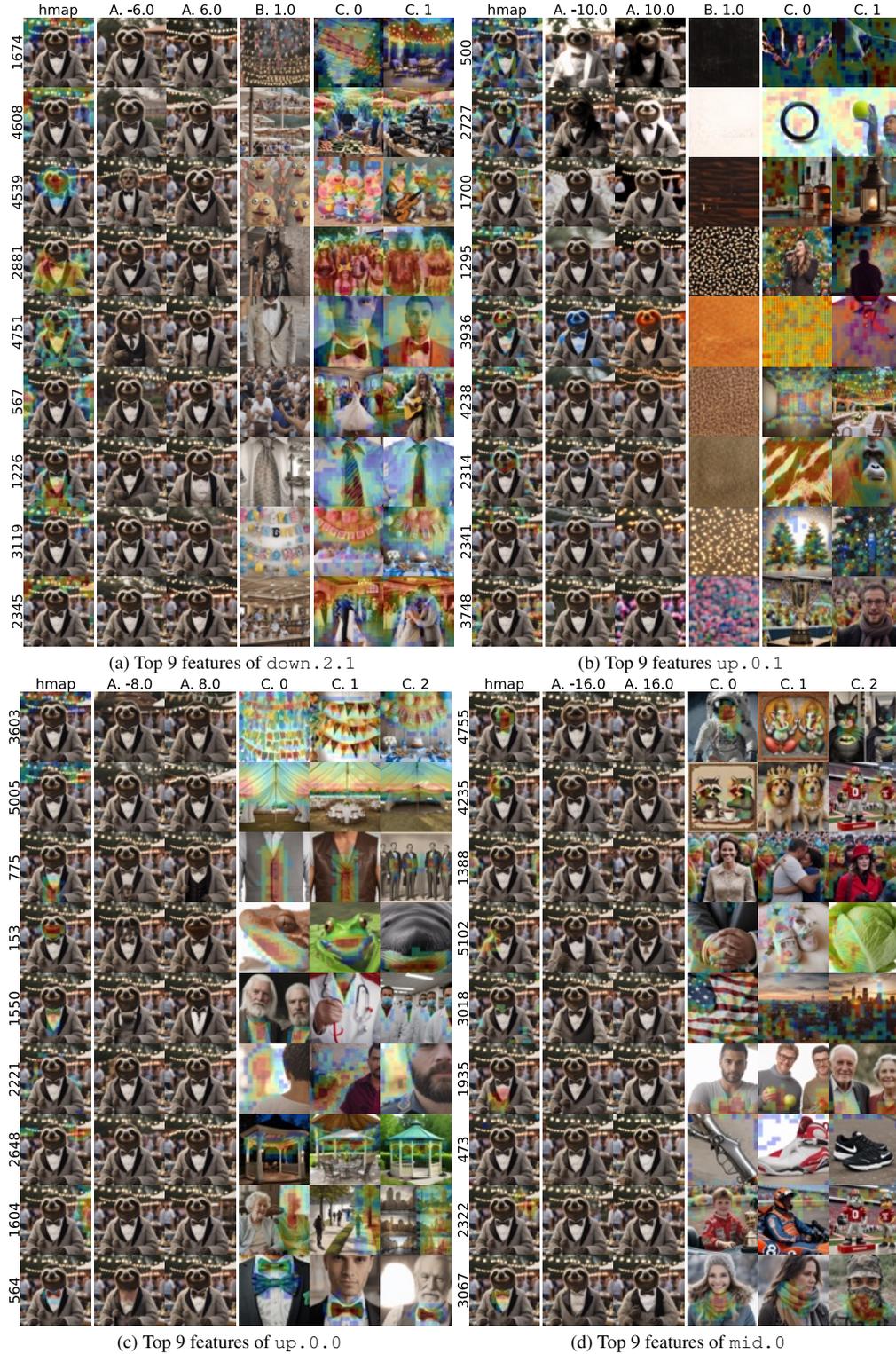


Figure 10. The top 9 features of `down.2.1` (a), `up.0.1` (b), `up.0.0` (c) and `mid.0` (d) for the prompt: "A cinematic shot of a professor sloth wearing a tuxedo at a BBQ party." Each row represents a feature. The first column depicts a feature heatmap (highest activation red and lowest nonzero one blue). The column titles containing "A" show feature modulation interventions, the ones containing "B" the intervention of turning on the feature on the empty prompt, and the ones containing "C" depict top dataset examples. Floating point values in the title denote β and γ values.

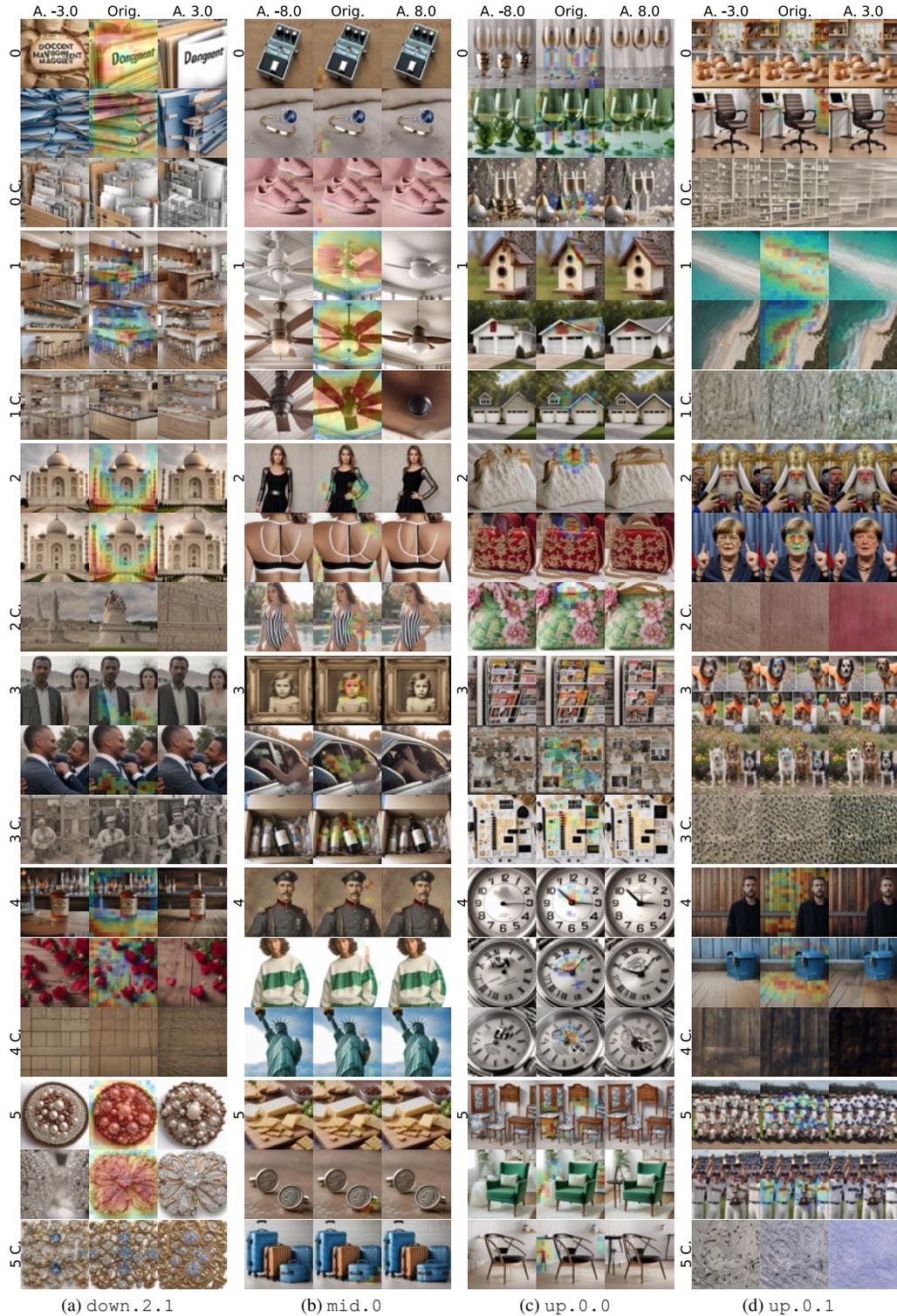


Figure 11. We visualize 6 features for down.2.1 (a), mid.0 (b), up.0.0, and up.0.1. We use three columns for each transformer block and three rows for each feature. For down.2.1 and up.0.1 we visualize the two samples from the top 5% quantile of activating dataset examples (middle) together a feature ablation (left) and a feature enhancement (right), and, activate the feature on the empty prompt with $\gamma = 0.5, 1, 2$ from left to right. For mid.0 and up.0.0 we display three samples with ablation and enhancement. Captions are in Table 5.

Table 5. Prompts for the top 5% quantile examples in Fig. 11

Block	Feature	Prompt
down.2.1	0	A file folder with the word document management on it.
	0	Two blue folders filled with dividers.
	1	A kitchen with an island and bar stools.
	1	An unfinished bar with stools and a wood counter.
	2	The Taj Mahal, or a white marble building in India.
	2	The Taj Mahal, or a white marble building in India.
	3	A man and woman standing next to each other.
	3	Two men in suits hugging each other outside.
	4	An old Forester whiskey bottle sitting on top of a wooden table.
	4	Red roses and hearts on a wooden table.
	5	A beaded brooch with pearls and copper.
	5	An image of a brooch with diamonds.
	mid.0	0
0		An engagement ring with blue sapphire and diamonds.
0		The women's pink sneaker is shown.
1		A white ceiling fan with three blades.
1		A ceiling fan with three blades and a light.
1		The ceiling fan is dark brown and has two wooden blades.
2		The black dress is made from knit and has metallic sleeves.
2		The back view of a woman wearing a black and white sports bra.
2		The woman is wearing a striped swimsuit.
3		An old-fashioned photo frame with a little girl on it.
3		The woman is sitting in her car with her head down.
3		The contents of an empty bottle in a box.
4		An old painting of a man in uniform.
4		The model wears an off-white sweatshirt with green panel.
4		The Statue of Liberty stands tall in front of a blue sky.
5	Cheese and crackers on a cutting board.	
5	Two cufflinks with coins on them.	
5	Three pieces of luggage are shown in blue.	
up.0.0	0	Three wine glasses with gold and silver designs.
	0	Three green wine glasses sitting next to each other.
	0	New Year's Eve with champagne, gold, and silver.
	1	The birdhouse is made from wood and has a brown roof.
	1	The garage is white with red shutters.
	1	Two garages with one attached porch and the other on either side.
	2	An elegant white lace purse with gold clasp.
	2	The red handbag has gold and silver designs.
	2	A pink and green floral-colored purse.
	3	A magazine rack with magazines on it.
	3	The year-in-review page for this digital scrap.
	3	The planner sticker kit is shown with gold and black accessories.
	4	A clock with numbers on the face.
	4	A silver watch with roman numerals on the face.
	4	An automatic watch with a silver dial.
5	Four pieces of wooden furniture with blue and white designs.	
5	The green chair is in front of a white rug.	
5	The wish chair with a black seat.	
up.0.1	0	The wooden toy kitchen set includes bread, eggs, and flour.
	0	The office chair is brown and black.
	1	An aerial view of the white sand and turquoise water.
	1	An aerial view of the beach and ocean.
	2	The patriarch of Ukraine is shown speaking to reporters.
	2	German Chancellor Merkel gestures as she speaks to the media.
	3	Four pictures showing dogs wearing orange vests.
	3	Two dogs are standing on the ground next to flowers.
	4	A man standing in front of a wooden wall.
	4	A blue mailbox sitting on top of a wooden floor.
5	The baseball players are posing for a team photo.	
5	The baseball players are holding up their trophies.	

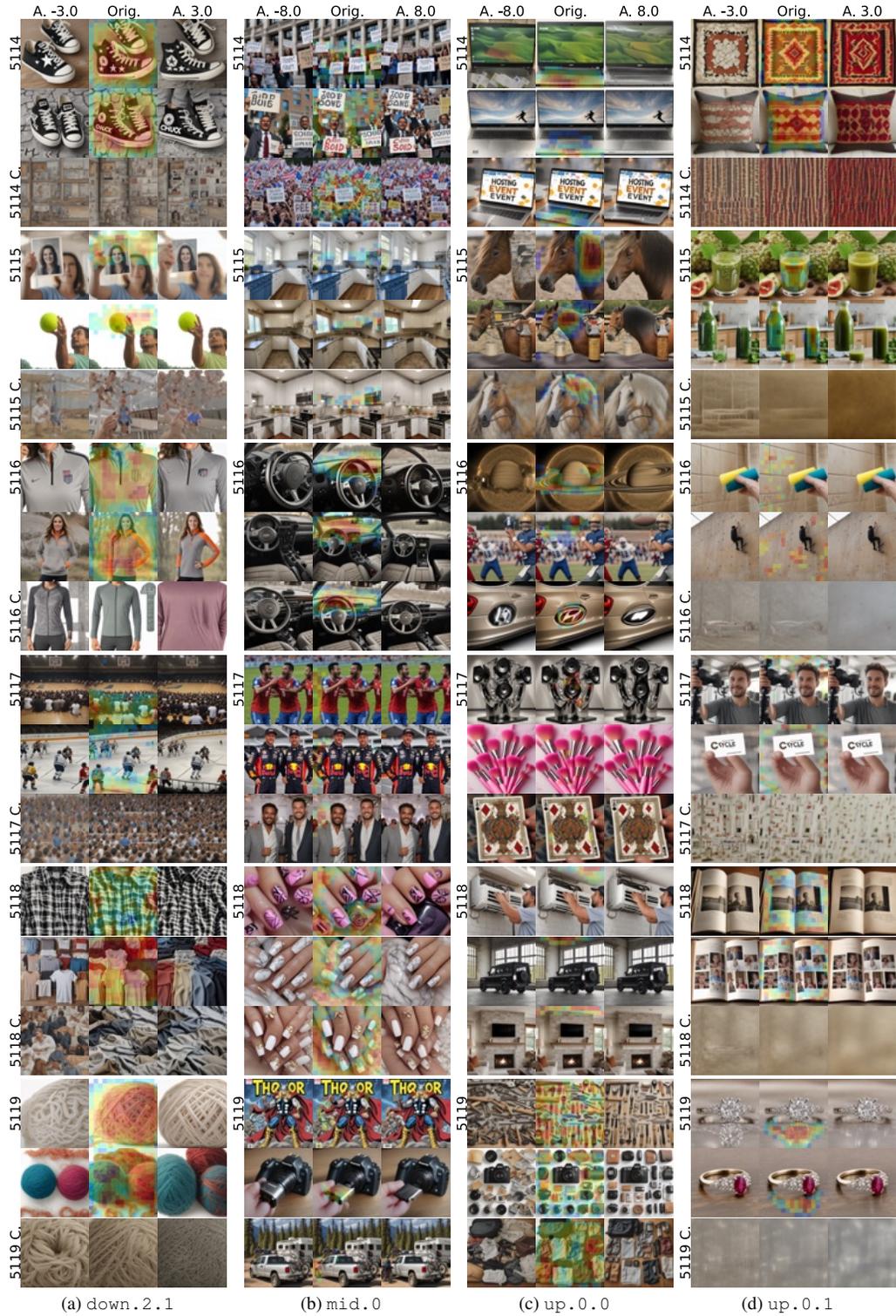


Figure 12. We visualize last 6 features for down.2.1 (a), mid.0 (b), up.0.0, and up.0.1. We use three columns for each transformer block and three rows for each feature. For down.2.1 and up.0.1 we visualize two samples from the top 5% quantile of activating dataset examples (middle) together a feature ablation (left) and a feature enhancement (right), and, activate the feature on the empty prompt with $\gamma = 0.5, 1, 2$ from left to right. For mid.0 and up.0.0 we display three samples with ablation and enhancement. Captions are in Table 6.

Table 6. Prompts for the top 5 % quantile examples in Fig. 12

Block	Feature	Prompt	
down.2.1	5114	Black and white Converse sneakers with the word black star.	
	5114	Black and white Converse sneakers with the word Chuck.	
	5115	A woman holding up a photo of herself.	
	5115	A man holding up a tennis ball in the air.	
	5116	The Nike Women’s U.S. Soccer Team DRI-Fit 1/4 Zip Top.	
	5116	The women’s gray and orange half-zip sweatshirt.	
	5117	A large group of people sitting in front of a basketball court.	
	5117	Hockey players are playing in an arena with spectators.	
	5118	The black and white plaid shirt is shown.	
	5118	The different colors and sizes of t-shirts.	
	5119	A ball of yarn on a white background.	
	5119	Two balls of colored wool are on the white surface.	
	mid.0	5114	People holding signs in front of a building.
		5114	Two men dressed in suits and ties are holding up signs.
5114		A large group of people holding flags and signs.	
5115		A kitchen with white cabinets and a blue stove.	
5115		The kitchen is clean and ready for us to use.	
5115		A kitchen with white cabinets and stainless steel appliances.	
5116		The steering wheel and dashboard in a car.	
5116		The interior of a car with dashboard controls.	
5116		The dashboard and steering wheel in a car.	
5117		Three men are celebrating a goal on the field.	
5117		Two men in Red Bull racing gear standing next to each other.	
5117		Two men are posing for the camera at an event.	
5118		Someone is holding up their nail polish with pink and black designs.	
5118		The nail is very cute and looks great with marble.	
5118	White stily nails with gold and diamonds.		
5119	The Mighty Thor comic book.		
5119	The camera is showing its flash drive.		
5119	A truck with bikes on the back parked next to a camper.		
up.0.0	5114	The Acer laptop is open and ready to use.	
	5114	The Lenovo S13 laptop is open and has an image of a person jumping off the keyboard.	
	5114	A laptop with the words Hosting Event on it.	
	5115	A horse with a black nose and brown mane.	
	5115	The horse leather oil is being used to protect horses.	
	5115	An oil painting on a canvas of a horse.	
	5116	The sun is shining brightly over Saturn.	
	5116	A football player throws the ball to another team.	
	5116	Car door light logo sticker for Hyundai.	
	5117	An artistic black and silver sculpture with speakers.	
	5117	The pink brushes are sitting on top of each other.	
	5117	Four kings playing cards in the hand.	
	5118	A man is fixing an air conditioner.	
	5118	The black Land Rover is parked in front of a large window.	
5118	A flat screen TV mounted on the wall above a fireplace.		
5119	A table with many different tools on it.		
5119	A camera with many different items including flash cards, lenses, and other accessories.		
5119	The contents of an open suitcase and some clothes.		
up.0.1	5114	An old Navajo rug with multicolored designs.	
	5114	The pillow is made from an old kilim.	
	5115	An image of noni juice with some fruits.	
	5115	A bottle and glass on the counter with green juice.	
	5116	Someone cleaning the shower with a sponge.	
	5116	A man on a skateboard climbing a wall with ropes.	
	5117	A man taking a selfie in front of some camera equipment.	
	5117	A person holding up a business card with the words cycle transportation.	
	5118	Two photos are placed on top of an open book.	
	5118	An open book with pictures of children and their parents.	
	5119	An engagement ring with diamonds on top.	
5119	An oval ruby and diamond ring.		

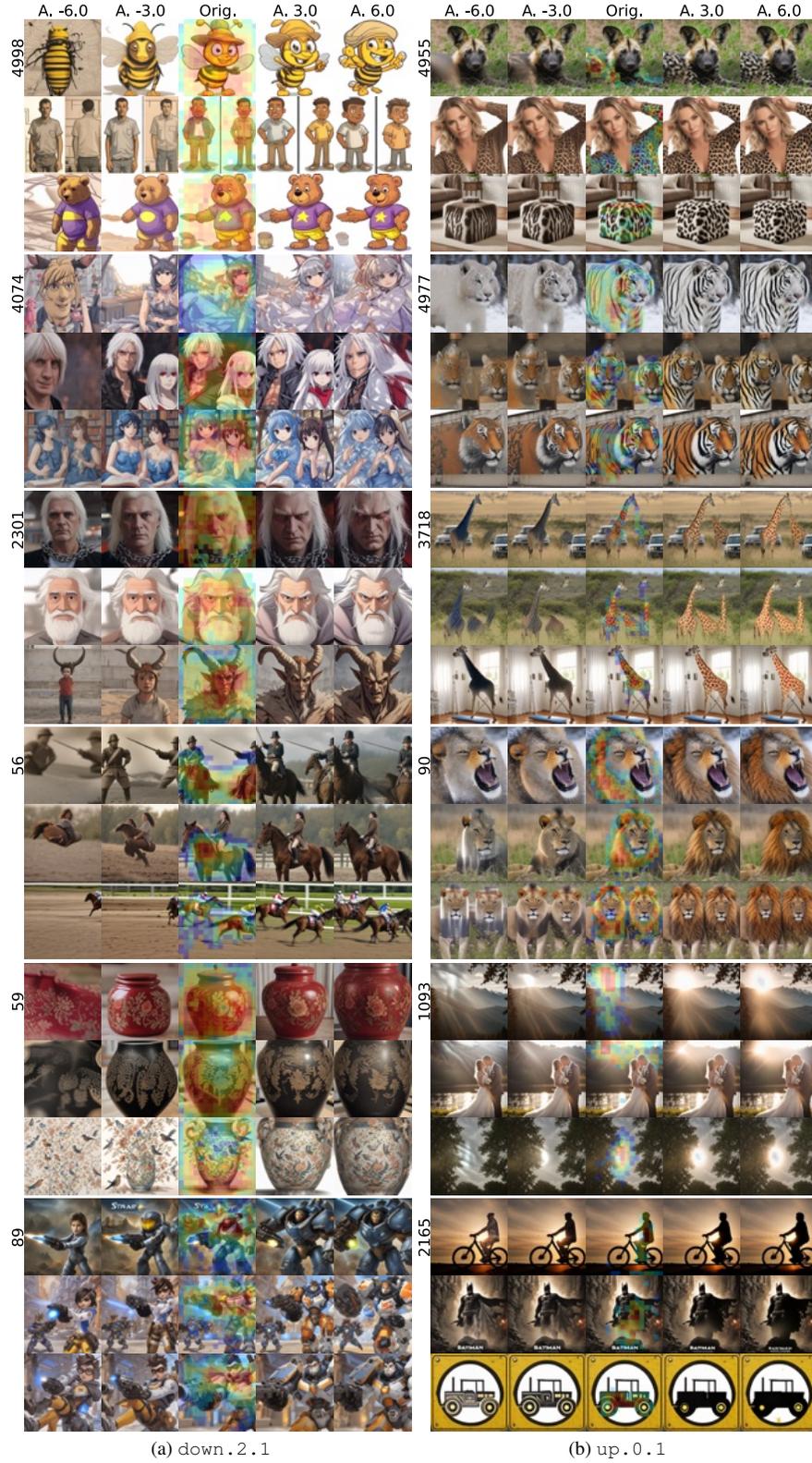


Figure 13. We visualize 6 features for down.2.1 (a) and up.0.1 (b). For each feature, we use 5 columns showing ablations (left), activating examples (middle), enhancements (right) and 3 rows with different samples from the top 5% quantile of activating examples. Captions are in Table 7.

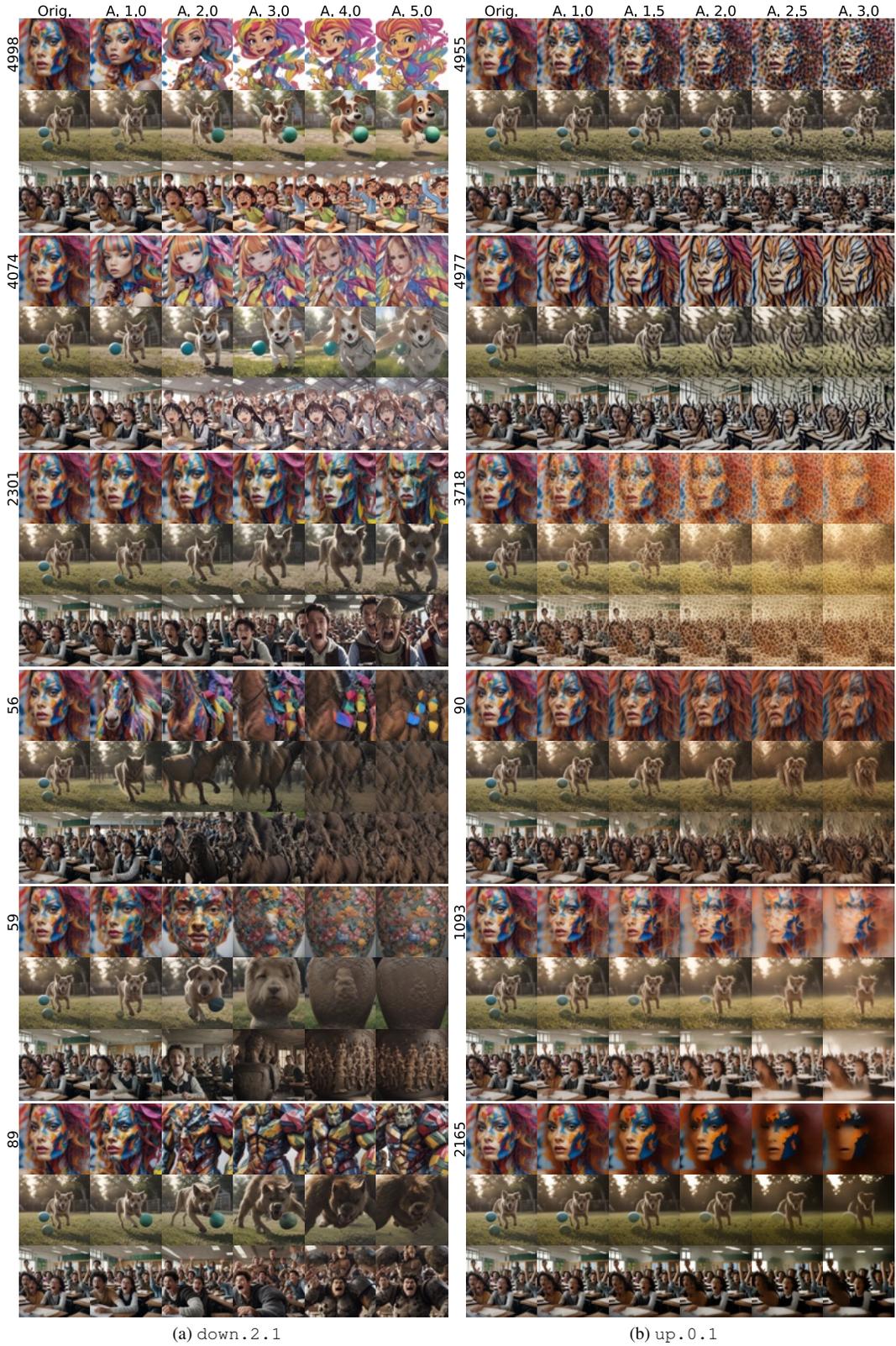


Figure 14. We turn on the features from Fig. 13 on three unrelated prompts “a photo of a colorful model”, “a cinematic shot of a dog playing with a ball”, and “a cinematic shot of a classroom with excited students”.

Table 7. Prompts for the top 5% quantile examples in Fig. 13

Block	Feature	Prompt	
down.2.1	4998	A cartoon bee wearing a hat and holding something.	
	4998	Two cartoon pictures of the same man with his hands in his pockets.	
	4998	A cartoon bear with a purple shirt and yellow shorts.	
	4074	An anime character with cat ears and a dress.	
	4074	Two anime characters, one with white hair and the other with red eyes.	
	4074	An anime book with two women in blue dresses.	
	2301	A man with white hair and red eyes holding a chain.	
	2301	An animated man with white hair and a beard.	
	2301	The character is standing with horns on his head.	
	56	Two men in uniforms riding horses with swords.	
	56	A woman riding on the back of a brown horse.	
	56	Two jockeys on horses racing down the track.	
	59	A red jar with floral designs on it.	
	59	An old black vase with some design on it.	
	59	A vase with birds and flowers on it.	
	89	StarCraft 2 is coming to the Nintendo Wii.	
	89	Overwatch is coming to Xbox and PS3.	
	89	The hero in Overwatch is holding his weapon.	
	up.0.1	4955	An African wild dog laying in the grass.
		4955	The woman is posing for a photo in her leopard print top.
4955		An animal print cube ottoman with brown and white fur.	
4977		A white tiger with blue eyes standing in the snow.	
4977		A bottle and tiger are shown next to each other.	
4977		A mural on the side of a building with a tiger.	
3718		Giraffes are standing in the grass near a vehicle.	
3718		Two giraffes standing next to each other in the grass.	
3718		A giraffe standing next to an ironing board.	
90		A lion is roaring its teeth in the snow.	
90		A lion sitting in the grass looking off into the distance.	
90		Two lions with flowers on their backs.	
1093		The sun is shining over mountains and trees.	
1093		Bride and groom in front of a lake with sun flare.	
1093		The milky sun is shining brightly over the trees.	
2165		The silhouette of a person riding a bike at sunset.	
2165	The Dark Knight rises from his cave in Batman's poster.		
2165	A yellow sign with black design depicting a tractor.		

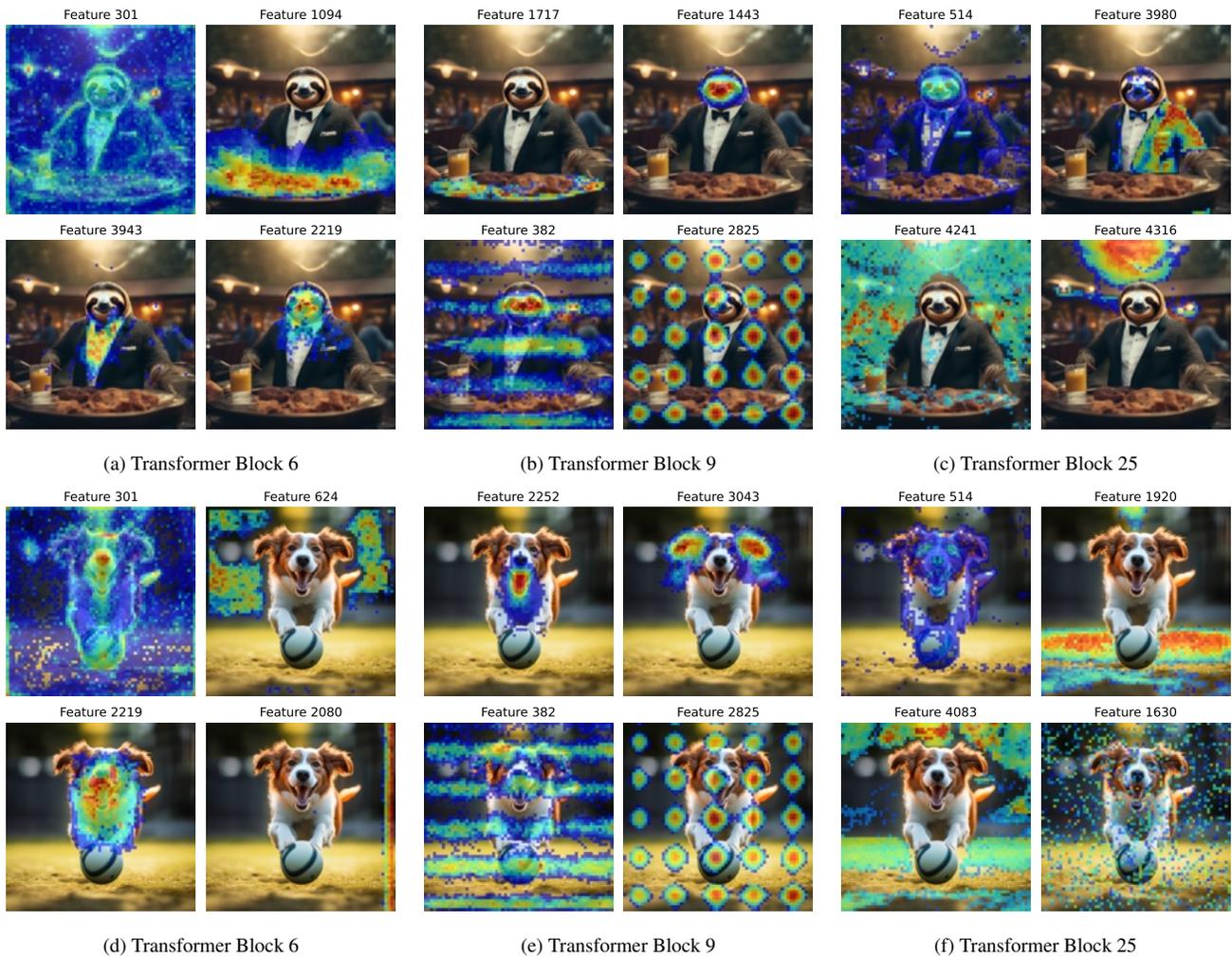


Figure 15. Activation maps of SAEs trained on Pixart-LCM XL transformer blocks. The top row corresponds to the prompt: "A cinematic shot of a sloth wearing a tuxedo at a BBQ party." The bottom row corresponds to the prompt: "A dog playing with a ball."