

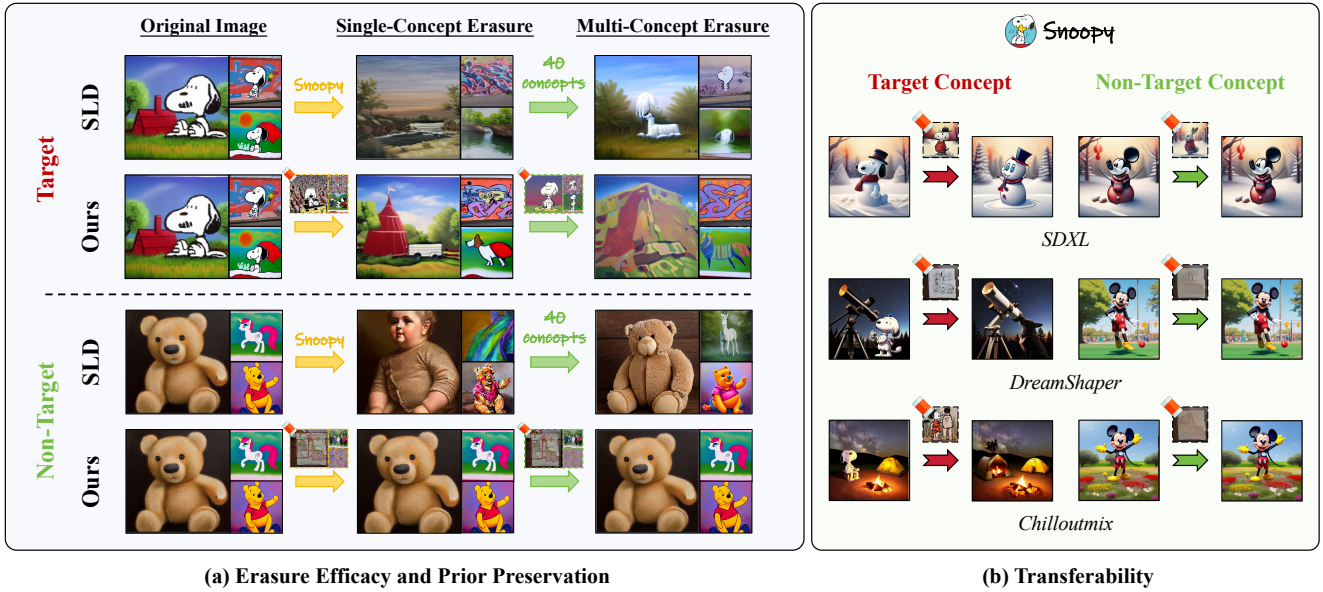
# Precise, Fast, and Low-cost Concept Erasure in Value Space: Orthogonal Complement Matters

Yuan Wang<sup>1\*</sup>, Ouxiang Li<sup>1\*</sup>, Tingting Mu<sup>2</sup>, Yanbin Hao<sup>3†</sup>, Kuien Liu<sup>4</sup>, Xiang Wang<sup>1</sup>, Xiangnan He<sup>1†</sup>

<sup>1</sup>University of Science and Technology of China, <sup>2</sup>The University of Manchester,

<sup>3</sup>Hefei University of Technology, <sup>4</sup>Institute of Software Chinese Academy of Sciences

{wy1001, lioox}@mail.ustc.edu.cn, haoyanbin@hotmail.com, xiangnanhe@gmail.com



(a) Erasure Efficacy and Prior Preservation

(b) Transferability

Figure 1. The proposed Adaptive Value Decomposer (AdaVD) demonstrates a satisfactory balance between erasure efficacy and prior preservation and an effective transferability across T2I diffusion models. (a) Compared to SLD [39], AdaVD enables precise concept erasure without compromising prior knowledge for non-target concepts at both single- and multi-concept erasure. This is facilitated by a precise disentanglement of target semantics (e.g., “Snoopy”) and a robust preservation of non-target ones (e.g., “Teddy”), with visualization interpretation marked by  $\diamond$ . (b) AdaVD can be transferred to various T2I models, e.g., SDXL [31], DreamShaper [6], Chilloutmix [5].

## Abstract

Recent success of text-to-image (T2I) generation and its increasing practical applications, enabled by diffusion models, require urgent consideration of erasing unwanted concepts, e.g., copyrighted, offensive, and unsafe ones, from the pre-trained models in a precise, timely, and low-cost manner. The twofold demand of concept erasure includes not only a precise removal of the target concept (i.e., erasure efficacy) but also a minimal change on non-target content (i.e., prior preservation), during generation. Existing methods face challenges in maintaining an effective balance between erasure efficacy and prior preservation, and they can be computationally costly. To improve, we propose a

precise, fast, and low-cost concept erasure method, called **Adaptive Value Decomposer (AdaVD)**, which is training-free. Our method is grounded in a classical linear algebraic operation of computing the orthogonal complement, implemented in the value space of each cross-attention layer within the UNet of diffusion models. We design a shift factor to adaptively navigate the erasure strength, enhancing effective prior preservation without sacrificing erasure efficacy. Extensive comparative experiments with both training-based and training-free state-of-the-art methods demonstrate that the proposed AdaVD excels in both single and multiple concept erasure, showing 2 to 10 times improvement in prior preservation than the second best, meanwhile achieving the best or near best erasure efficacy. AdaVD supports a series of diffusion models and downstream image generation tasks, with code available on:

\*Equal Contributions.

†Corresponding authors.

## 1. Introduction

The recent advancements of text-to-image (T2I) diffusion models [11, 17, 18, 30, 36, 50, 52] have enabled users to effortlessly generate high-quality images with simple textual prompts. However, such generations would inevitably introduce copyrighted [10, 22, 42, 44] or offensive [23, 39] concepts, caused by the noisy training data scraped from web [9, 40]. Because it is very costly to re-train large generative models from scratch, it is vital to develop low-cost techniques to precisely erase unwanted semantic concepts in images, *i.e.*, *concept erasure*. This task aims at a precise erasure of visual content w.r.t. target concepts from generated images (*i.e.*, *erasure efficacy*), while a faithful preservation of irrelevant content w.r.t. the prompts comprising non-target concepts (*i.e.*, *prior preservation*), safeguarding a secure T2I generation for diffusion models.

A representative category of concept erasure methods is training-based, which fine-tunes a subset of model parameters by formulating meticulous erasing objective functions [12, 20, 26, 27]. Despite good erasure efficacy, they exhibit considerable practical limitations when being deployed. For instance, they require expensive individual fine-tuning to erase each concept, which thereby limits their real-time usage. As an example, it is unacceptable for online T2I platforms to be costly to erase newly emerging concepts, where copyrighted or offensive concepts could arise unexpectedly, with no means to produce a complete list of concepts to erase in advance. Moreover, these methods suffer from a limited balance between erasure efficacy and prior preservation, due to their reliance on regularization terms to trade off prior preservation.

An alternative category of concept erasure methods is training-free, such as Negative Prompt (NP) [2], Safe Latent Diffusion (SLD) [39], and SuppressEOT [24], which enable real-time erasure. They intervene in the image generation process, exhibiting a range of drawbacks. For instance, NP was initially designed to enhance image quality and can result in compromised erasure efficacy; while SuppressEOT requires the user to specify the location of the target concept within the prompt, thus, it is not suitable for erasure applications that require full automation. Therefore, both NP and SuppressEOT fall short as independent tools for concept erasure. Regarding SLD, it does not perform well at retaining prior knowledge of non-target concepts as illustrated in Fig. 1, failing in *precise concept erasure*. Limited by their drawbacks, current training-free methods are not robust enough to act as an independent concept erasure tool and to be applied in continual erasure of multiple concepts.

In this light, we advance concept erasure techniques for T2I generation, by developing a precise, fast, and low-cost method called **Adaptive Value Decomposer (AdaVD)**. It is

a training-free method, capable of precisely erasing the target concepts and satisfactorily preserving non-target priors with low computational overhead. Our core design builds on a classical linear algebraic operation, *i.e.*, projection onto the orthogonal complement of the semantic space of the target concepts. We conduct this projection-based decomposition in the cross-attention value space, disentangling target semantics from the original prompts. To improve the erasure precision and prior preservation, we further refine the decomposition by adaptively allocating token-wise shifts. These shifts are designed to differentiate the strong and specific alignments from the weak and general alignments between prompt tokens and visual content associated with the target concept. Guided by these shifts, the strong and specific alignments are erased for precision while the weak and general alignments are retained for prior preservation.

Fig. 1 (a) demonstrates the erased results for both target and non-target concepts. It shows that our AdaVD can precisely locate and erase the components indicated by the target semantics, meanwhile keeping the non-target priors maximally unaffected. Empirical evaluation shows that AdaVD excels in prior preservation, outperforming the second best by a 2- to 10-fold improvement across various non-target concepts, and meanwhile maintains exceptional erasure efficacy, consistently achieving the best or near-best performance. Our contributions are summarized below:

- A novel and effective erasing operation by exploiting the projection onto the orthogonal complement of the target concept in the cross-attention value space, to disentangle semantics carried by the target concept.
- An adaptive erasing mechanism through a dynamic shift factor, which can effectively minimize the impact on prior knowledge, without compromising the erasure efficacy.
- A precise, fast, and low-cost concept erasure technique AdaVD, which works in a training-free manner and supports a series of T2I diffusion models, capable of precise concept erasure.
- Extensive experiments that demonstrate the superiority of AdaVA against state-of-the-art (SOTA) methods, achieving 2 to 10 times of improvement in prior preservation while maintaining precise erasure efficacy.

## 2. Related Works

**Re-training and Blocking:** The most straightforward way to erase a target concept from a pre-trained T2I model is to exclude the training data relevant to this concept and re-train the model from scratch, as in Stable Diffusion (SD) v2.0 [3]. For this, an Not Safe For Work (NSFW) detector [8, 21, 38] can be used to filter unsafe data from LAION-5B [41] prior to the training. However, this approach is time-consuming, requires specialized detectors, and can introduce biases [43]. An alternative solution is to block the prompts of concerns and restrict the outputs of concerns,



by using filters [43] and safety checkers [1, 35]. Yet, such safeguards are fragile and easy to bypass [46, 49], especially when being confronted with crafted malicious prompts [39].

**Training-based:** A more effective group of concept erasure solutions fine-tune pre-trained generative models, teaching them to “forget” a target concept, e.g., Erased Stable Diffusion (ESD) [12]. However, since the ESD design does not consider prior preservation, both target and non-target concepts are adversely affected. To improve, new training techniques have been developed, e.g., ConAbl [20], SA [14], UCE [13], EraseDiff [48], SPM [27] and MACE [26]. They improve prior preservation through regularization but still fall short in achieving both precise erasure and robust prior preservation. Moreover, for every new target concept to be erased, a separate fine-tuning is required. This is time-consuming, making real-time erasure impractical, especially for highly interactive platforms where unsafe or inappropriate concepts can emerge unexpectedly.

**Training-free:** With largely reduced computing costs, training-free methods are gaining increasing attention. NP [2] and SLD [39] pioneer training-free concept erasure by adjusting classifier-free guidance [17], leading the generation towards a direction away from the target concepts. However, NP lacks fine-grained control over target concepts and compromises prior preservation. SLD also compromises the prior knowledge during its generation, disrupting the overall generating quality of non-target concepts. SuppressEOT [24], akin to image editing techniques [15], removes target concepts based on user-specified textual positions. Its user-involved design makes it more suitable for editing tasks, but not for system-wide erasing tasks that require full automation. Benefiting from orthogonal decompositions in value spaces, the proposed training-free method AdaVD achieves not only precise concept erasure but also satisfactory prior preservation. It supports multi-concept erasure, is compatible across different versions of stable diffusion, and consistently demonstrates superior performance. The training-free method SAFREE [51], concurrent to ours, also uses orthogonal decomposition but operates differently on text embeddings, supported by masking, projection, Fourier transforms, and a hard control of removal strength. We compare performance with it in the appendix.

### 3. Method

Concept erasure in T2I generation aims at a successful removal of the visual content indicated by textual concepts (*i.e.*, target concepts), meanwhile a satisfactory preservation of the visual content irrelevant to these concepts (*i.e.*, prior knowledge). It is challenging to achieve simultaneously satisfactory erasure efficacy and prior preservation. To advance this field, we propose a precise, fast, and low-cost concept erasure method, termed AdaVD and illustrated in Fig. 2 (c). The method is training-free, and its core de-

sign builds on a classical linear algebraic operation, *i.e.*, orthogonal complement. This simple but elegant geometric operation is able to guide effectively the image generation, enabling a precise removal of the target concepts from the image content meanwhile a satisfactory preservation of the non-target content.

#### 3.1. Preliminary on T2I Diffusion Models

Current T2I models usually include an image compression network [19] and a conditional latent diffusion model [36] that performs sequential denoising with a UNet [37] in the latent space. The UNet takes as inputs the noise latent variable  $\mathbf{z}_t$ , timestep  $t$ , and embedding  $\mathbf{C}$  of the textual prompt from a pre-trained CLIP model [33], and predicts the noise  $\epsilon_\theta(\mathbf{z}_t, t, \mathbf{C})$ . In training-free concept erasure, the noise is additionally conditioned on the target concept with text embedding  $\mathbf{C}_t$ , predicted by  $\epsilon_\theta(\mathbf{z}_t, t, \mathbf{C}, \mathbf{C}_t)$ .

Interactions between the image and text modalities are enabled by cross-attention (CA) layers [36, 47] within the UNet, which align the latent representation of the noisy image with the semantic detail of the textual prompt. Each CA layer computes an attention map  $\mathbf{A} = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d}}\right)$  with a latent feature dimension  $d$ . The queries  $\mathbf{Q}$  are computed from the noisy image features, while both keys  $\mathbf{K}$  and values  $\mathbf{V}$  from the text embeddings  $\mathbf{C}$ , using different projection matrices. The layer output is a weighted aggregation of  $\mathbf{A}$  and  $\mathbf{V}$ . More details on CA layers are in Appendix A.

It has been recognized that the keys in CA layers mostly act as the “Where” pathway, governing the layout of the attention map and determining the compositional structure of the generated images, while the values as the “What” pathway, controlling the content and visual appearance of images [45]. Because the goal of concept erasure is to modify the visual content of the generated images, we propose to conduct value decompositions into subspaces, which are uniquely constructed by exploiting the target concept and its orthogonal complement in an adaptive fashion. We demonstrate that information offered by the orthogonal complement can be used to generate successfully high-quality images with the target concept precisely erased.

#### 3.2. Token-wise Target Embedding Pre-processing

Given a textual example, which can be either a target concept to erase or an original prompt, its embedding is computed at the token level by a CLIP text encoder. Each tokenized example is padded with [SOT] as the prefix and [EOT] at the end, with [EOT] filling any remaining positions to maintain a fixed token length of  $l$ . Each token is characterized by an embedding vector of dimension  $D_c$ .

We focus on the embedding of a target concept  $\mathbf{C}_t \in \mathbb{R}^{l \times D_c}$ , and denote each column of  $\mathbf{C}_t^T$  by  $\mathbf{c}_t^j$  that corresponds to a token vector where  $j \in \{1, 2, \dots, l\}$ . To facilitate a precise erasure, we emphasize the key informa-

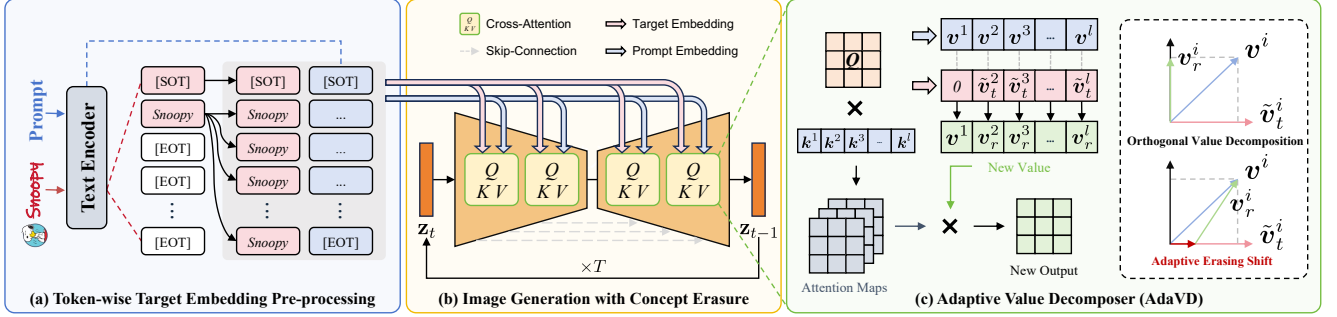


Figure 2. **Overview of our Adaptive Value Decomposer (AdaVD)** in erasing the target concept “Snoopy”. (a) First, we token-wisely duplicate the last subject token of the target embedding encoded by the text encoder, except for [SOT]. (b) Then, the pre-processed target embedding and corresponding prompt embedding are jointly fed into CA layers within the UNet as conditions, to disentangle target semantics from the original image at each timestep. (c) In each CA layer, we perform token-wise orthogonal value decomposition with an adaptive token-wise shift. The new value is subsequently multiplied by the attention map, producing the erased output for this CA layer.

tion carried by a target concept by proposing an embedding duplication. Specifically, the embedding of the last subject token within its prompt content excluding [SOT] and [EOT], replaces all the token positions except for [SOT], as illustrated in Fig. 2 (a). For instance, the single-token concept “snoopy” has a modified embedding matrix corresponding to “[SOT], snoopy, snoopy, ..., snoopy”, while the multi-token concept “Van Gogh” to “[SOT], gogh, gogh, ..., gogh”. Benefiting from the causal attention mechanism used by the CLIP text encoder, the last subject token is able to “see” all the prompt content and contains key information [28], therefore is sufficient for erasure calculation. The modified embedding matrix  $\tilde{\mathbf{C}}_t$  is fed into a CA layer to compute the value matrix  $\tilde{\mathbf{V}}_t \in \mathbb{R}^{l \times d}$  via a linear projection using the projection matrix  $\mathbf{W}_V \in \mathbb{R}^{D_e \times d}$ , as

$$\tilde{\mathbf{V}}_t = \tilde{\mathbf{C}}_t \mathbf{W}_V = \begin{bmatrix} \mathbf{c}_t^1, \underbrace{\mathbf{c}_t^k, \dots, \mathbf{c}_t^k}_{l-1} \end{bmatrix}^T \mathbf{W}_V. \quad (1)$$

The token vector  $\mathbf{c}_t^1$  corresponds to [SOT] and  $\mathbf{c}_t^k$  to the last subject token. The total token number is still  $l$ .

### 3.3. Orthogonal Value Decomposition

Given a conditional latent diffusion model trained for standard T2I generation, our proposed erasing operation works by projecting the original textual prompt onto the orthogonal complement of the subspace spanned by the target concepts to erase, and it is implemented in the value space learned at each CA layer of the UNet. It supports both single-concept and multi-concept erasure.

We do not apply any erasing operation to [SOT], as it primarily serves as a prefix and does not carry useful information to distinguish the semantic content. We start from the modified value matrix  $\tilde{\mathbf{V}}_t$  in Eq. (1), and use  $\tilde{\mathbf{v}}_t^j$  to denote the  $j$ -th column of  $\tilde{\mathbf{V}}_t^T$ . The exclusion of [SOT] in the erasure calculation is equivalent to replacing the value

vector in the first row of  $\tilde{\mathbf{V}}_t$  by a zero vector, resulting in  $\mathbf{V}_t = [\mathbf{0}, \tilde{\mathbf{v}}_t^2, \tilde{\mathbf{v}}_t^3, \dots, \tilde{\mathbf{v}}_t^l]^T$  to use in the erasing operation.

#### 3.3.1. Single-concept Erasure

Our erasing operation works with  $\mathbf{V}_t^T$ , where each column of  $\mathbf{V}_t^T$  corresponds to the erased value vector for token position  $j$  and is denoted as  $\mathbf{v}_t^j$ , and with the value matrix  $\mathbf{V} \in \mathbb{R}^{l \times d}$  computed from the original prompt, where each column of  $\mathbf{V}^T$  is referred to as the original value vector for the token position  $j$ , denoted by  $\mathbf{v}^j$ . To remove the effect of the target concept from the original prompt, for each token position, we project the original value vector  $\mathbf{v}^j$  onto the orthogonal complement of the span of the erased value vector  $\mathbf{v}_t^j$ , and denote this orthogonal complement by  $\text{span}^\perp(\mathbf{v}_t^j)$ . This results in the following modified value vector:

$$\begin{aligned} \mathbf{v}_r^j &= \mathbf{P}_{\text{span}^\perp(\mathbf{v}_t^j)} \mathbf{v}^j = (\mathbf{I}_d - \mathbf{P}_{\text{span}(\mathbf{v}_t^j)}) \mathbf{v}^j \\ &= \mathbf{v}^j - \frac{\mathbf{v}_t^j \mathbf{v}_t^{jT}}{\mathbf{v}_t^{jT} \mathbf{v}_t^j} \mathbf{v}^j = \mathbf{v}^j - \frac{\mathbf{v}_t^{jT} \mathbf{v}^j}{\mathbf{v}_t^{jT} \mathbf{v}_t^j} \mathbf{v}_t^j, \end{aligned} \quad (2)$$

where  $\mathbf{P}_{\mathbb{X}} \mathbf{x}$  denotes the orthogonal projection of a vector  $\mathbf{x}$  onto the space  $\mathbb{X}$ , and  $\mathbf{I}_d$  is an identity matrix of size  $d$ . The modified value vector  $\mathbf{v}_r^j$  is used, instead of  $\mathbf{v}^j$ , to calculate the output of the CA layer, as illustrated in Fig. 2 (c). Since  $\mathbf{v}_t^1 = \mathbf{0}$ , it has  $\mathbf{v}_r^1 = \mathbf{v}^1$ , meaning that no erasure is performed for [SOT].

#### 3.3.2. Multi-concept Erasure

We generalize the above operation to erase a set of  $n$  target concepts, and their corresponding modified value matrices are denoted by  $\{\mathbf{V}_t^h \in \mathbb{R}^{l \times d}\}_{h=1}^n$ . We use  $\mathbf{v}_t^{h,j}$  to denote the  $j$ -th column of  $(\mathbf{V}_t^h)^T$ , referred to as the erased value vector for the  $j$ -th token position of the  $h$ -th target concept. Our erasing operation can naturally be extended to projecting the original prompt to the orthogonal complement of the span of the  $n$  erased value vectors  $\{\mathbf{v}_t^{h,j}\}_{h=1}^n$ ,

denoted by  $\text{span}^\perp \left( \left\{ \mathbf{v}_t^{h,j} \right\}_{h=1}^n \right)$ . To calculate the projection, we first conduct the Gram-Schmidt orthogonalization to obtain a set of  $n$  orthonormal basis vectors  $\left\{ \mathbf{o}_t^{h,j} \right\}_{h=1}^n$  for the span of  $\left\{ \mathbf{v}_t^{h,j} \right\}_{h=1}^n$ . Here, we assume the value vectors in  $\left\{ \mathbf{v}_t^{h,j} \right\}_{h=1}^n$  are linearly independent and they form a basis. Such an assumption is reasonable because the multiple concepts to erase should be semantically different in practice, otherwise, a single-concept erasure would be sufficient. The desired projection is then computed by

$$\begin{aligned} \mathbf{v}_r^j &= \mathbf{P}_{\text{span}^\perp(\{\mathbf{v}_t^{h,j}\}_{h=1}^n)} \mathbf{v}^j = \mathbf{P}_{\text{span}^\perp(\{\mathbf{o}_t^{h,j}\}_{h=1}^n)} \mathbf{v}^j \quad (3) \\ &= \left( \mathbf{I}_d - \mathbf{P}_{\text{span}(\{\mathbf{o}_t^{h,j}\}_{h=1}^n)} \right) \mathbf{v}^j \\ &= \mathbf{v}^j - \sum_{h=1}^n \left( \mathbf{o}_t^{h,j} \right)^T \mathbf{v}^j \mathbf{o}_t^{h,j}. \end{aligned}$$

As an addition, we provide in Appendix B an alternative way to compute  $\mathbf{P}_{\text{span}^\perp(\{\mathbf{v}_t^{h,j}\}_{h=1}^n)} \mathbf{v}^j$  that does not require Gram-Schmidt orthogonalization but matrix inverse.

### 3.4. Adaptive Erasing Shift

In practice, given a pair of textual prompts and a target concept to erase, their token-wise relevance can vary across different token positions. The different tokens of the prompt carry information with different intensities and focuses and thus can have quite different effects on image generation. As an example, we compare the generated images in each subfigure of Fig. 3, using different versions of the value matrix  $\mathbf{V}$  for the same prompt. We separate the values corresponding to the [EOT] tokens and the remaining, expressing it as  $\mathbf{V} = [\mathbf{V}_{\text{content}}, \mathbf{V}_{\text{[EOT]}}]$ . Herein, we include [SOT] within  $\mathbf{V}_{\text{content}}$  for simplicity. We obtain three versions of  $\mathbf{V}$  by (1) keeping it as what it is, (2) setting  $\mathbf{V}_{\text{[EOT]}}$  as zero, and (3) setting  $\mathbf{V}_{\text{content}}$  as zero. Fig. 3 shows that the prompt content carries more featured information than those [EOT] tokens. This motivates us to further improve the design in Eqs. (2) and (3), by enabling adaptive adjustment of the erasing operation at the token level.

We discover that, although a projection onto the orthogonal complement of the target concept is effective at erasing this concept itself, it can sometimes affect the prior preservation due to an excessive removal of information. Therefore, our adaptive design is focused on improving the prior preservation. When semantics carried by a prompt token are less relevant to a target token, we intend to perform less erasure to protect the prior image content. According to [32], angular information plays a more critical role in conveying semantics than magnitude during image generation. These motivate us to exploit the cosine similarity between the value vectors of a prompt token and a target token and use it to derive a shift factor to adjust the erasing strength

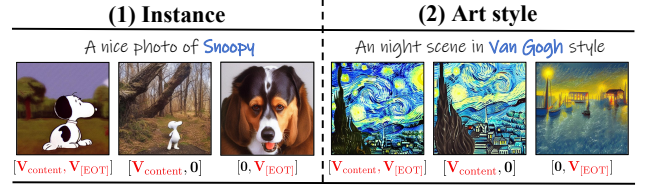


Figure 3. We analyze the contribution of different tokens in text-visual alignment by separately masking the value of content tokens and [EOT] tokens, where content tokens carry more featured information than those [EOT] tokens.

along the identified erasing direction, which we refer to as the erasing shift. In general, we let the factor reduce when the cosine similarity becomes smaller.

We denote the shift factor by  $\delta(\cdot, \cdot) : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ , which is a function of two value vectors. Using it, we revise the erasing operation in Eq. (2) for a single concept to

$$\mathbf{v}_r^j = \mathbf{v}^j - \frac{\delta(\mathbf{v}_t^j, \mathbf{v}^j) \mathbf{v}_t^{jT} \mathbf{v}^j}{\mathbf{v}_t^{jT} \mathbf{v}_t^j} \mathbf{v}_t^j. \quad (4)$$

While for multiple concepts, Eq. (3) is revised to

$$\mathbf{v}_r^j = \mathbf{v}^j - \sum_{h=1}^n \delta(\mathbf{v}_t^{h,j}, \mathbf{v}^j) \left( \sum_{k=1}^n w_{hk} \left( \mathbf{o}_t^{k,j} \right)^T \mathbf{v}^j \right) \mathbf{v}_t^{h,j}, \quad (5)$$

where  $w_{hk}$  is the  $hk$ -th element of the projection matrix that transforms the basis  $\left\{ \mathbf{v}_t^{h,j} \right\}_{h=1}^n$  to the orthonormal basis  $\left\{ \mathbf{o}_t^{h,j} \right\}_{h=1}^n$ . For a single target concept, it is straightforward to introduce the shift factor, as in Eq. (4). However, in the case of erasing multiple concepts, the derivation of a shifted operation is less straightforward, as the orthonormal basis does not carry meaningful semantics, thus a transformation back to the value vectors of the target tokens is needed. We explain in Appendix B how to derive Eq. (5) from Eq. (3). The revised erasing operations result in the proposal AdaVD.

Our shift factor design builds on the sigmoid function. Given two input vectors  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ , we formulate it as

$$\delta(\mathbf{x}, \mathbf{y}) = \frac{s}{1 + e^{-p(\cos(\mathbf{x}, \mathbf{y}) - \epsilon)}}. \quad (6)$$

The cosine threshold  $\epsilon$  allows a strong erasure when exceeding the threshold. A negative cosine similarity indicates a very weak relevance between the target and prompt tokens, suggesting no erasure is required. Therefore, we set  $0 < \epsilon < 1$  to quantify and filter the relatively weak relevance indicated by a positive but small cosine similarity. The hyper-parameter  $s > 0$  controls the factor scale, while  $p > 0$  controls the increasing rate of the  $\delta$  value against the cosine value. In Appendix C.2, we provide hyper-parameter implementation details and a comprehensive analysis.



## 4. Experiments and Result Analysis

We conduct extensive experiments to evaluate the proposed AdaVD for erasing a diverse range of target concepts, covering specific instances, art styles, NSFW content, and celebrity. We compare with SOTA training-based methods, including ConAbl [20], ESD [12], SPM [27], MACE [26] and SOTA training-free methods including NP [2] and SLD [39]. We also demonstrate the time efficiency and interpretability of AdaVD, along with its wider usage in other downstream image generation tasks, coupled with a series of diffusion models.

### 4.1. Experimental Setup

**Implementation:** We employ SD v1.4 [2] to generate images using the DPM-solver sampler [25] over 30 sampling steps with classifier-free guidance [17] of 7.5. All the compared methods are implemented following their default configurations available from their official repository. Further implementation details are provided in Appendix C.1.

**Evaluation Data:** Adopting the same evaluation protocol from SPM [27], we assess the methods based on 80 instance templates, 30 art style templates, and 25 celebrity templates, and benchmark each method by generating 10 images per template per concept in evaluation. To assess the performance of NSFW erasure, we use the I2P benchmark [39].

**Performance Metrics:** We follow the widely used evaluation metrics for concept erasure, including the CLIP score (CS) [33] to assess erasure efficacy and the Fréchet inception distance (FID) [16] to assess prior preservation. CS calculates the cosine similarity between a textual prompt and the generated image [33]. We examine two CS values before and after erasing and compare the CS value after the erasure. When the prompts contain the target concepts, a more reduced CS indicates a more effective erasure of the target concepts. FID measures the distance between images generated before and after erasing [16]. A lower FID indicates a better alignment between the two images. Thus, for non-target concepts, lower FID values indicate better prior preservation. Overall, a precise concept erasure should have a low CS for prompts containing the target concepts and a low FID for prompts composed of non-target concepts.

**Summary:** Results on instance and art style concept erasure are reported in Sections 4.2 and 4.3, while results on celebrity and NSFW erasure are reported in Appendices D.3 and D.2. Additional analyses, including extended quantitative results, comparisons with SuppressEOT [24], performance on different versions of SD, and further discussions, are presented in the Appendix.

### 4.2. On Instance Concept Erasure

We first conduct the experiment with erasing a single concept “Snoopy”. Six types of prompts were tested, of which one contains “Snoopy” and the other five contain only non-

Concept	Snoopy	Mickey	Spongebob	Pikachu	Dog	Legislator
	CS	CS	CS	CS	CS	CS
SD v1.4	28.51	26.57	27.43	-	-	-
<i>Erase Snoopy</i>						
	CS ↓	FID ↓	FID ↓	FID ↓	FID ↓	FID ↓
ConAbl	25.38	38.44	41.59	29.68	27.76	27.36
MACE	<u>20.78</u>	118.01	111.90	81.99	43.27	65.97
SPM	23.89	<u>33.06</u>	<u>34.70</u>	<u>23.89</u>	<u>19.61</u>	<u>18.26</u>
NP	23.66	59.58	78.74	52.37	67.51	55.22
SLD	27.84	48.12	55.36	38.74	41.95	49.08
Ours	<b>20.28</b>	<b>5.72</b>	<b>8.56</b>	<b>5.79</b>	<b>2.32</b>	<b>6.07</b>
<i>Erase Snoopy and Mickey</i>						
	CS ↓	CS ↓	FID ↓	FID ↓	FID ↓	FID ↓
ConAbl	24.26	24.08	46.32	39.63	30.57	27.49
MACE	<u>20.74</u>	<u>20.71</u>	51.49	110.67	52.07	77.13
SPM	23.16	22.81	<u>41.58</u>	<u>31.77</u>	<u>21.96</u>	<u>23.69</u>
NP	23.59	24.85	81.41	50.10	65.93	58.88
SLD	27.76	26.74	54.59	39.24	41.62	50.13
Ours	<b>20.29</b>	<b>19.93</b>	<b>9.34</b>	<b>5.84</b>	<b>2.41</b>	<b>6.43</b>
<i>Erase Snoopy and Mickey and Spongebob</i>						
	CS ↓	CS ↓	CS ↓	FID ↓	FID ↓	FID ↓
ConAbl	23.94	23.64	25.04	51.20	31.59	30.03
MACE	<u>20.48</u>	<u>20.50</u>	21.59	99.68	47.46	70.38
SPM	22.81	22.35	<u>20.82</u>	39.83	<u>22.68</u>	<u>25.31</u>
NP	24.29	24.76	25.31	64.75	65.10	59.33
SLD	27.84	26.71	27.60	<u>39.41</u>	42.32	49.88
Ours	<b>19.39</b>	<b>19.73</b>	<b>20.34</b>	<b>6.85</b>	<b>2.79</b>	<b>7.26</b>

Table 1. **Quantitative comparison of single- and multi-instance erasure.** The best and second-best results are marked in **bold** and underlined, respectively. Our AdaVD consistently achieves the lowest CS and the lowest FID in all cases, indicating superior prior preservation without compromising erasure efficacy.

target concepts. The results are compared in the top block of Table 1. It can be seen that our proposed AdaVD achieves the lowest CS and FID scores in all cases. Particularly, its FID is 33% lower than that of the second-best method. It improves over SOTA by significantly enhanced prior preservation without compromising the erasure precision, as exemplified in Fig. 4. It can be observed from Fig. 4 that methods such as SPM, NP, and SLD fail to fully erase “Snoopy” which can be found from the ear and the shape characteristic of the generated image after erasure. MACE can successfully erase the “Snoopy” concept. But all the competing methods suffer from degraded image quality in non-target concept generation, e.g., “Spongebob”.

We then compare the performance for multi-concept erasure, experimenting with two cases, of which one erases two concepts of “Snoopy” and “Mickey” together and the other erases three concepts of “Snoopy”, “Mickey” and “Spongebob” together. Results are reported in the bottom two blocks of Table 1 and visualized in the second and third sections of Fig. 4. Similarly, our AdaVD achieves the lowest CS and FID scores across all cases. This verifies that AdaVD is capable of handling more complex multi-concept erasure and simultaneously fighting against catastrophic forgetting. Conversely, other methods struggle to



Figure 4. **Qualitative comparison of single- and multi-instance erasure.** Both training-based and training-free methods show limitations in prior preservation. In contrast, our AdaVD demonstrates considerable performance in maintaining prior knowledge without compromising erasure efficacy across both single- and multi-concept erasure tasks.

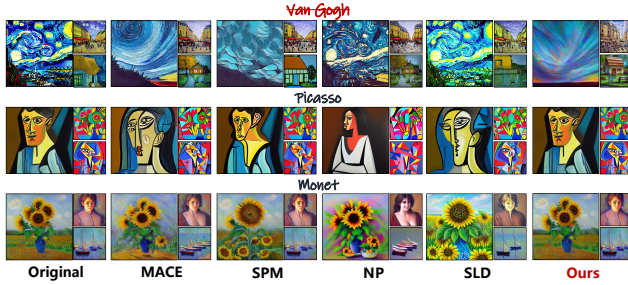


Figure 5. **Qualitative comparison of art style erasure.** Our AdaVD can effectively remove the target concept “Van Gogh” while preserving non-target styles like “Picasso” and “Monet”.

consistently perform well for multi-concept erasure. For example, when erasing “Snoopy” and “Mickey” together, the image quality for non-target concepts like “Spongebob” and “Legislator” shows apparent degradation in Fig. 4. SLD, in particular, fails to maintain its erasure ability even when dealing with 2-concept erasure. In Appendix G, we report more experimental results for erasing up to 40 concepts, where our AdaVD can extend to erase dozens of concepts in practice, maintaining consistent erasure efficacy and prior preservation, as shown in Fig. 1.

### 4.3. On Art Style Erasure

We experiment with erasing specific art style, including “Van Gogh”, “Picasso” and “Monet”. Results are reported in Table 2, and visual comparisons are provided in Fig. 5. Our AdaVD exhibits superior prior preservation, and achieves the lowest or close-to-lowest CS and FID scores, demonstrating strong prior preservation without sacrificing erasure efficacy. Other methods show different drawbacks as observed from the results. For instance, although NP achieves notably better precision in art style removal as

Concept	Van Gogh	Picasso	Monet	Andy Warhol	Caravaggio
	CS	CS	CS	CS	CS
SD v1.4	29.21	29.06	29.02	-	-
<i>Erase Van Gogh</i>					
	CS ↓	FID ↓	FID ↓	FID ↓	FID ↓
ConAbl	28.80	71.71	138.72	70.30	73.10
MACE	27.74	65.77	69.79	83.37	75.41
SPM	<b>24.78</b>	62.25	32.27	58.30	61.50
NP	24.90	141.56	124.52	127.85	136.32
SLD	27.48	103.96	109.11	103.89	119.32
Ours	<u>24.87</u>	<b>6.82</b>	<b>2.66</b>	<b>8.36</b>	<b>6.84</b>
<i>Erase Picasso</i>					
	FID ↓	CS ↓	FID ↓	FID ↓	FID ↓
ConAbl	58.62	27.72	140.34	73.35	67.44
MACE	60.46	27.11	49.92	76.10	72.85
SPM	<u>38.79</u>	<u>26.69</u>	<u>7.76</u>	<u>52.00</u>	<u>51.40</u>
NP	111.35	<b>26.14</b>	91.11	116.24	121.82
SLD	98.21	27.03	93.01	97.00	110.05
Ours	<b>5.49</b>	26.99	<b>2.33</b>	<b>9.38</b>	<b>7.05</b>
<i>Erase Monet</i>					
	FID ↓	FID ↓	CS ↓	FID ↓	FID ↓
ConAbl	141.52	132.10	<u>24.53</u>	208.38	186.26
MACE	76.90	69.35	26.89	88.35	81.72
SPM	<u>41.03</u>	<u>29.71</u>	27.00	<u>31.90</u>	<u>25.99</u>
NP	137.21	126.75	<b>24.47</b>	127.22	135.83
SLD	94.48	92.88	25.73	100.90	114.87
Ours	<b>6.94</b>	<b>6.50</b>	26.30	<b>8.46</b>	<b>7.19</b>

Table 2. **Quantitative comparison of art style erasure.** AdaVD achieves a superior balance between erasure efficacy and prior preservation, especially excelling in prior preservation.

compared to instance removal, Fig. 5 shows that it still struggles to fully erase the “Van Gogh” style. Also, SLD fails to erase the “Van Gogh” style. These two methods also degrade generation quality for non-target styles, such as “Picasso” and “Monet”, showing harmful effects on non-target concepts directly, as evidenced in both Fig. 5 and

	Data Preparation	Model Finetune	Image Generation	Total Time
ConAbl	9290	1120	0.9	10419
SPM	0	72850	1.7	72867
MACE	303	232	0.9	544
SLD	0	0	1.4	14
Ours	4	0	1.8	22

Table 3. **Time consumption of 10-concept erasure.** We calculate the time cost (s) to erase 10 concepts and generate 10 images using one NVIDIA A40 GPU. Compared with training-based methods, AdaVD exhibits exceptional efficiency in real-time erasure.

Table 2. Both MACE and SPM are effective in erasing the target concept, however, their prior preservation is somehow less satisfactory, which is particularly noticeable in the generated images in “*Monet*” style. Differently, AdaVD can effectively and consistently remove the target concept and meanwhile preserve satisfactorily the prior content, as confirmed by Table 2 and Fig. 5.

#### 4.4. Further Analysis

**Time Consumption:** The computational cost of concept erasure primarily arises from three components, including (1) *data preparation time* required by training-based methods for preparing training data and by AdaVD for basis computation; (2) *model fine-tuning time* required by training-based methods; and (3) *image generation time* required by all methods. In Table 3, we compare the total time consumption of different methods, as well as their time spent on each component. The two training-free methods of SLD and AdaVD are significantly faster as no fine-tuning is needed. Our AdaVD costs slightly more time than SLD, *i.e.*, 0.8 extra seconds per image due to its basis computation, and a total of 8 extra seconds for generating 10 images. But this mild increase yields a significant performance gain, succeeding in precise concept erasure.

**Interpreting Erased Components by Visualization:** Our AdaVD generates images by replacing the original value vector  $v^j$  with  $v_r^j$  via orthogonal complement operation. To empirically interpret the rationale behind our method, we visualize the erased component,  $v^j - v_r^j$ . Fig. 6 presents three cases of erasing the target concepts “*Mickey*”, “*Van Gogh*” and “*Bruce Lee*”, where AdaVD successfully erases these target concepts while robustly preserving prior knowledge of non-target concepts. In the first row, as shown in the middle column of each block, the erased components consistently align with the corresponding target semantics when dealing with the target concepts. In the second row of non-target concepts, conversely, the erased components do not contain any informative pattern, indicating that they carry no meaningful semantics, and therefore exert minimal impact on the prior knowledge.

**Downstream Applications:** We conduct additional experiments to showcase the versatile applications of our AdaVD across various generative tasks, including (1) implicit con-



Figure 6. **Visualization of erased components.** In each block, we compare both target (1st row) and non-target concept (2nd row) by visualizing the original image (1st column), erased component (2nd column), and generation by our AdaVD (3rd column).

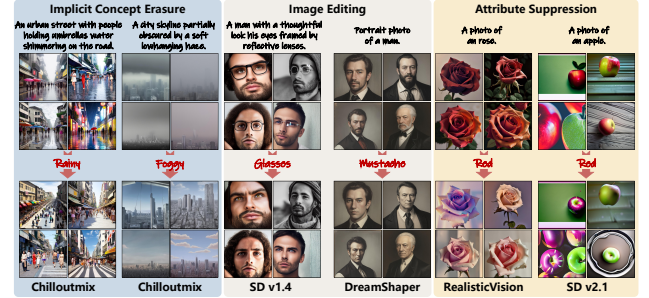


Figure 7. **Downstream applications.** We extend AdaVD to versatile generation tasks, including (1) implicit concept erasure, (2) image editing, and (3) attribute suppression, indicating its significant potential for broad applications.

cept erasure: by removing the implicit concepts of “*rainy*” and “*foggy*”; (2) image editing: by removing the appearance concepts of “*glasses*” and “*mustache*”; and (3) attribute suppression: by removing the coupled color concept of “*red*”. Additionally, we integrate AdaVD with a series of diffusion models, including Chilloutmix [5], DreamShaper [6], RealisticVision [7], and SD v2.1 [4], alongside SD v1.4. As illustrated in Fig. 7, despite the absence of explicit mention for “*rainy*” and “*foggy*”, AdaVD can still effectively erase these concepts in image semantic space. Meanwhile, AdaVD also precisely removes “*glasses*” and “*mustache*” with minimal changes to other details, highlighting its potential in image editing applications. For attribute suppression, AdaVD successfully eliminates the color attribute “*red*” from objects such as apples and roses, demonstrating its capability to decouple strongly coupled concepts, *e.g.*, “*roses are red*” embedded in the model’s prior knowledge.

## 5. Conclusion and Future Work

We have presented AdaVD, a precise, fast, and low-cost method for erasing unwanted concepts. The idea of leveraging the classical linear algebraic orthogonal complement operation and an adaptive erasing shift design is novel, and has successfully achieved a precise concept erasure. Extensive experiments have demonstrated both high erasure efficacy and strong prior preservation of AdaVD for both single- and multi-concept erasure. Moreover, AdaVD exhibits excellent interpretability through visualizing its



erased components and strong capability in solving downstream tasks. It has been an intriguing discovery that the orthogonal complement of the value vectors of the target concepts is effective at erasing their inherent semantics. Despite the empirical success, we will seek to establish a rigorous theoretical understanding of the accumulative effect of applying this linear algebraic operation layer-wise for concept erasure in the future.

## Acknowledgements

This work is supported by the National Science and Technology Major Project (2023ZD0121102) and the National Natural Science Foundation of China (U24B20180)

## References

- [1] Safety checker nested in stable diffusion. <https://huggingface.co/CompVis/stable-diffusion-safety-checker>, 2022. 3
- [2] Stable diffusion. <https://huggingface.co/CompVis/stable-diffusion-v1-4>, 2022. 2, 3, 6
- [3] Stable diffusion v2.0. <https://huggingface.co/stabilityai/stable-diffusion-2>, 2022. 2
- [4] Stable diffusion v2.1. <https://huggingface.co/stabilityai/stable-diffusion-2-1>, 2022. 8
- [5] Chilloutmix. <https://huggingface.co/swl-models/chilloutmix>, 2023. 1, 8, 19
- [6] Dreamshaper. <https://huggingface.co/Lykon/DreamShaper>, 2023. 1, 8, 19, 22
- [7] Realisticvision. [https://huggingface.co/SG161222/Realistic\\_Vision\\_V5.1\\_noVAE](https://huggingface.co/SG161222/Realistic_Vision_V5.1_noVAE), 2023. 8, 19, 22
- [8] Praneeth Bedapudi. Nudenet: Neural nets for nudity detection and censoring, 2022. URL <https://github.com/notAI-tech/NudeNet>. 2
- [9] Abeba Birhane, Vinay Uday Prabhu, and Emmanuel Kahembwe. Multimodal datasets: misogyny, pornography, and malignant stereotypes. *arXiv preprint arXiv:2110.01963*, 2021. 2
- [10] Nicolas Carlini, Jamie Hayes, Milad Nasr, Matthew Jagielski, Vikash Sehwal, Florian Tramer, Borja Balle, Daphne Ippolito, and Eric Wallace. Extracting training data from diffusion models. In *32nd USENIX Security Symposium (USENIX Security 23)*, pages 5253–5270, 2023. 2
- [11] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021. 2
- [12] Rohit Gandikota, Joanna Materzynska, Jaden Fiotto-Kaufman, and David Bau. Erasing concepts from diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2426–2436, 2023. 2, 3, 6
- [13] Rohit Gandikota, Hadas Orgad, Yonatan Belinkov, Joanna Materzynska, and David Bau. Unified concept editing in diffusion models. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 5111–5120, 2024. 3
- [14] Alvin Heng and Harold Soh. Selective amnesia: A continual learning approach to forgetting in deep generative models. *Advances in Neural Information Processing Systems*, 36, 2024. 3
- [15] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-or. Prompt-to-prompt image editing with cross-attention control. In *The Eleventh International Conference on Learning Representations*. 3
- [16] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017. 6
- [17] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. In *NeurIPS 2021 Workshop on Deep Generative Models and Downstream Applications*. 2, 3, 6
- [18] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 2
- [19] Diederik P Kingma. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. 3
- [20] Nupur Kumari, Bingliang Zhang, Sheng-Yu Wang, Eli Shechtman, Richard Zhang, and Jun-Yan Zhu. Ablating concepts in text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 22691–22702, 2023. 2, 3, 6
- [21] Gant Laborde. Deep nn for nsfw detection. 2
- [22] Ouxiang Li, Jiayin Cai, Yanbin Hao, Xiaolong Jiang, Yao Hu, and Fuli Feng. Improving synthetic image detection towards generalization: An image transformation perspective. *arXiv preprint arXiv:2408.06741*, 2024. 2
- [23] Ouxiang Li, Yanbin Hao, Zhicai Wang, Bin Zhu, Shuo Wang, Zaixi Zhang, and Fuli Feng. Model inversion attacks through target-specific conditional diffusion models. *arXiv preprint arXiv:2407.11424*, 2024. 2
- [24] Senmao Li, Joost van de Weijer, Fahad Khan, Qibin Hou, Yaxing Wang, et al. Get what you want, not what you don’t: Image content suppression for text-to-image diffusion models. In *The Twelfth International Conference on Learning Representations*. 2, 3, 6, 21
- [25] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver: A fast ode solver for diffusion probabilistic model sampling in around 10 steps. *Advances in Neural Information Processing Systems*, 35:5775–5787, 2022. 6
- [26] Shilin Lu, Zilan Wang, Leyang Li, Yanzhu Liu, and Adams Wai-Kin Kong. Mace: Mass concept erasure in diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6430–6440, 2024. 2, 3, 6
- [27] Mengyao Lyu, Yuhong Yang, Haiwen Hong, Hui Chen, Xuan Jin, Yuan He, Hui Xue, Jungong Han, and Guiguang Ding. One-dimensional adapter to rule them all: Concepts diffusion models and erasing applications. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7559–7568, 2024. 2, 3, 6, 22, 23
- [28] Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual associations in gpt.

- Advances in Neural Information Processing Systems*, 35: 17359–17372, 2022. 4
- [29] Carl D Meyer. *Matrix analysis and applied linear algebra*. SIAM, 2023. 11
- [30] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *International conference on machine learning*, pages 8162–8171. PMLR, 2021. 2
- [31] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. In *The Twelfth International Conference on Learning Representations*. 1, 17
- [32] Zeju Qiu, Weiyang Liu, Haiwen Feng, Yuxuan Xue, Yao Feng, Zhen Liu, Dan Zhang, Adrian Weller, and Bernhard Schölkopf. Controlling text-to-image diffusion by orthogonal finetuning. *Advances in Neural Information Processing Systems*, 36:79320–79362, 2023. 5
- [33] 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. 3, 6
- [34] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67, 2020. 18
- [35] Javier Rando, Daniel Paleka, David Lindner, Lennart Heim, and Florian Tramèr. Red-teaming the stable diffusion safety filter. *arXiv preprint arXiv:2210.04610*, 2022. 3
- [36] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 2, 3
- [37] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical image computing and computer-assisted intervention—MICCAI 2015: 18th international conference, Munich, Germany, October 5-9, 2015, proceedings, part III 18*, pages 234–241. Springer, 2015. 3
- [38] Patrick Schramowski, Christopher Tauchmann, and Kristian Kersting. Can machines help us answering question 16 in datasheets, and in turn reflecting on inappropriate content? In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*, pages 1350–1361, 2022. 2
- [39] Patrick Schramowski, Manuel Brack, Björn Deiseroth, and Kristian Kersting. Safe latent diffusion: Mitigating inappropriate degeneration in diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22522–22531, 2023. 1, 2, 3, 6
- [40] Christoph Schuhmann, Robert Kaczmarczyk, Aran Komatsuzaki, Aarush Katta, Richard Vencu, Romain Beaumont, Jenia Jitsev, Theo Coombes, and Clayton Mullis. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. In *NeurIPS Workshop Datacentric AI*, number FZJ-2022-00923. Jülich Supercomputing Center, 2021. 2
- [41] 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. 2
- [42] Shawn Shan, Jenna Cryan, Emily Wenger, Haitao Zheng, Rana Hanocka, and Ben Y Zhao. Glaze: Protecting artists from style mimicry by {Text-to-Image} models. In *32nd USENIX Security Symposium (USENIX Security 23)*, pages 2187–2204, 2023. 2
- [43] Zhan Shi, Xu Zhou, Xipeng Qiu, and Xiaodan Zhu. Improving image captioning with better use of captions. *arXiv preprint arXiv:2006.11807*, 2020. 2, 3
- [44] Gowthami Somepalli, Vasu Singla, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Diffusion art or digital forgery? investigating data replication in diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6048–6058, 2023. 2
- [45] Yoad Tewel, Rinon Gal, Gal Chechik, and Yuval Atzmon. Key-locked rank one editing for text-to-image personalization. In *ACM SIGGRAPH 2023 Conference Proceedings*, pages 1–11, 2023. 3
- [46] Yu-Lin Tsai, Chia-Yi Hsu, Chulin Xie, Chih-Hsun Lin, Jia You Chen, Bo Li, Pin-Yu Chen, Chia-Mu Yu, and Chun-Ying Huang. Ring-a-bell! how reliable are concept removal methods for diffusion models? In *The Twelfth International Conference on Learning Representations*. 3
- [47] A Vaswani. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017. 3
- [48] Jing Wu, Trung Le, Munawar Hayat, and Mehrtash Harandi. Erasediff: Erasing data influence in diffusion models. *arXiv preprint arXiv:2401.05779*, 2024. 3
- [49] Yijun Yang, Ruiyuan Gao, Xiaosen Wang, Tsung-Yi Ho, Nan Xu, and Qiang Xu. Mma-diffusion: Multimodal attack on diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7737–7746, 2024. 3
- [50] Hu Ye, Jun Zhang, Sibor Liu, Xiao Han, and Wei Yang. Ip-adapt: Text compatible image prompt adapter for text-to-image diffusion models. *arXiv preprint arXiv:2308.06721*, 2023. 2
- [51] Jaehong Yoon, Shoubin Yu, Vaidehi Patil, Huaxiu Yao, and Mohit Bansal. Safree: Training-free and adaptive guard for safe text-to-image and video generation. In *The Thirteenth International Conference on Learning Representations*. 3, 20
- [52] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3836–3847, 2023. 2

## A. Extra Preliminary on CA Layers

The cross-attention (CA) layers in the conditional denoising UNet  $\epsilon_\theta(\mathbf{z}_t, t, \mathbf{C})$  align the latent representation of the noisy image with that of the textual prompt. The latent variable at time step  $t$  is denoted by  $\mathbf{z}_t \in \mathbb{R}^{D_z \times H \times W}$  with a spatial dimension  $H \times W$  and a channel dimension  $D_z$ , while the text embedding, i.e. the latent representation of the textual prompt, is denoted by  $\mathbf{C} \in \mathbb{R}^{l \times D_c}$ . At the  $i$ -th CA layer, the attention map is computed by

$$\mathbf{A}_i = \text{softmax} \left( \frac{\mathbf{Q}_i \mathbf{K}_i^\top}{\sqrt{d_i}} \right), \quad (7)$$

where  $d_i$  is the latent feature dimension. The queries  $\mathbf{Q}_i \in \mathbb{R}^{H_i W_i \times d_i}$  are obtained by projecting the latent features of the noisy image returned by the previous module, while both the keys  $\mathbf{K}_i \in \mathbb{R}^{l \times d_i}$  and values  $\mathbf{V}_i \in \mathbb{R}^{l \times d_i}$  are computed by projecting the text embedding but using different projection matrices. Finally, the output of this CA layer is computed from the attention map, and the values by  $\mathbf{z}_t^{i+1} = \phi(\mathbf{A}_i \mathbf{V}_i)$ , where a common choice of  $\phi(\cdot)$  is a multi-layer perceptron. The subsequent modules take  $\mathbf{z}_t^{i+1}$  for further processing. For the convenience of explaining, we do not distinguish the notation  $i$  between layers in the main text.

## B. Equation Derivation

### B.1. On Equation (5)

Working with the subspace constructed as the span of the vector set  $\{\mathbf{v}_t^{h,j}\}_{h=1}^n$ , we obtain a set of orthonormal basis  $\{\mathbf{o}_t^{h,j}\}_{h=1}^n$  through the Gram-Schmidt orthogonalization. When the value vectors  $\{\mathbf{v}_t^{h,j}\}_{h=1}^n$  are linearly independent, each orthonormal basis can be expressed as a linear combination of these vectors such that

$$\mathbf{o}_t^{h,j} = \sum_{k=1}^n w_{hk} \mathbf{v}_t^{k,j}, \quad (8)$$

where  $w_{hk}$  are the combination weights. We have explained the linear independence assumption on  $\{\mathbf{v}_t^{h,j}\}_{h=1}^n$  in Section 3.3.2. Incorporating Eq. (8) into Eq. (3) but replacing only the second  $\mathbf{o}_t^{h,j}$ , it results in the following revised calculation of the orthogonal complement:

$$\mathbf{v}_r^j = \mathbf{v}^j - \sum_{h=1}^n \left( \sum_{k=1}^n w_{hk} (\mathbf{o}_t^{k,j})^T \mathbf{v}^j \right) \mathbf{v}_t^{h,j}. \quad (9)$$

The importance of this revised equation lies in the fact that it computes a weighted sum of the value vectors when performing the erasing. This enables the application of the adaptive erasing shift mechanism based on the value vectors, for which we further revise the erasing operation as

$$\mathbf{v}_r^j = \mathbf{v}^j - \sum_{h=1}^n \delta(\mathbf{v}_t^{h,j}, \mathbf{v}^j) \left( \sum_{k=1}^n w_{hk} (\mathbf{o}_t^{k,j})^T \mathbf{v}^j \right) \mathbf{v}_t^{h,j}. \quad (10)$$

Storing the combination weights in the matrix  $\mathbf{W} = [w_{hk}] \in \mathbb{R}^{n \times n}$ , it acts as a projection matrix transforming the two vector sets by

$$[\mathbf{o}_t^{1,j} \quad \dots \quad \mathbf{o}_t^{n,j}] = [\mathbf{v}_t^{1,j} \quad \dots \quad \mathbf{v}_t^{n,j}] \mathbf{W}. \quad (11)$$

### B.2. Alternative Orthonormal Basis Calculation

Purely for the interest of readers, we point out an alternative way to calculate the orthonormal basis. Constructing a matrix  $\hat{\mathbf{V}}_t^j \in \mathbb{R}^{d \times n}$  by using  $\{\mathbf{v}_t^{h,j}\}_{h=1}^n$  as its columns, following Equation (5.13.6) of the linear algebra textbook [29], the projection of  $\mathbf{v}^j$  onto  $\text{span}^\perp \left( \{\mathbf{v}_t^{h,j}\}_{h=1}^n \right)$  can be directly computed from  $\hat{\mathbf{V}}_t^j$  by

$$\begin{aligned} \mathbf{v}_r^j &= \mathbf{P}_{\text{span}^\perp(\{\mathbf{v}_t^{h,j}\}_{h=1}^n)} \mathbf{v}^j \\ &= \left( \mathbf{I}_d - \hat{\mathbf{V}}_t^j \left( (\hat{\mathbf{V}}_t^j)^T \hat{\mathbf{V}}_t^j \right)^{-1} (\hat{\mathbf{V}}_t^j)^T \right) \mathbf{v}^j. \end{aligned} \quad (12)$$



Compared to Eq. (3), Eq. (12) does not require the Gram-Schmidt orthogonalization, but the inverse calculation. Defining  $\mathbf{P}_t^j = \hat{\mathbf{V}}_t^j \left( \left( \hat{\mathbf{V}}_t^j \right)^T \hat{\mathbf{V}}_t^j \right)^{-1} \left( \hat{\mathbf{V}}_t^j \right)^T$ , one potential way to enable token-wise adaptive erasing shift based on Eq. (12) is

$$\mathbf{v}_r^j = \left( \mathbf{I}_d - \text{Diag} \left[ \delta \left( \mathbf{v}_t^{h,j}, \mathbf{v}^j \right) \right] \mathbf{P}_t^j \right) \mathbf{v}^j, \quad (13)$$

where  $\text{Diag} \left[ \delta \left( \mathbf{v}_t^{h,j}, \mathbf{v}^j \right) \right]$  is a diagonal matrix with shift factors  $\left[ \delta \left( \mathbf{v}_t^{h,j}, \mathbf{v}^j \right) \right]$  as its diagonal elements. We leave the in-depth investigation of exploiting this operation in practice to our future work.

## C. Additional Experimental Details

### C.1. On Implementation

To implement SD v1.4, the DPM-solver is chosen as the sampler, with a total of 30 sampling timesteps and a classifier-free guidance scale set of 7.5. Notably, we set the unconditional prompt to null text, as the negative prompt serves as a training-free method that can be directly compared with our AdaVD. To ensure a fair comparison, particularly for prior preservation, we use the same random seed (seed 0) across all methods to generate images under identical conditions. For the specific instance, art style, and celebrity erasure, we simply fix the hyperparameters to  $p = 100$ ,  $\epsilon = 0.93$ , and  $s = 2$ . Our AdaVD performs consistently well with this unified hyper-parameter configuration.

### C.2. Additional Hyper-parameter Analysis

The hyper-parameters of the shift factor, including  $0 < \epsilon < 1$  and  $p, s > 0$ , are closely related to the cosine similarities between tokens of the target concepts and tokens of the prompt. When erasing instances, art styles, and celebrity concepts, we notice that certain non-target concepts contained by the prompt semantically correlate with the target concept, with fairly strong correlations. For example, the non-target concept “Mickey” exhibits a relatively large cosine similarity of 0.65 with the target concept “Snoopy”, as they both belong to the category of cartoon characters. This makes it a fine balance between an unaffected generation of these non-target concepts and a successful erasure of the target concept. To examine how the erasure strength impacts such a balance, we show in Fig. 8 different image examples generated by AdaVD under various hyperparameter settings, for the target concept “Snoopy” and non-target concept “Mickey”.

Overall, the factor scale  $s$  and the threshold  $\epsilon$  significantly impact the balance between the erasure efficacy and prior preservation. Specifically, it can be observed, from the top part of Fig. 8 (on target concept), that a reduction in  $\epsilon$  results in a greater deviation in the generated images as compared to the original, for content relevant to the target concept. This indicates an enhanced erasure efficacy. This effect is further amplified as  $s$  increases. When adopting the setting of  $s = 2$  and  $\epsilon = 0.6$ , the erasure becomes excessive. Conversely, for non-target concept generation, a lower threshold  $\epsilon$  can negatively impact the non-target prior, as observed from the bottom part of Fig. 8 (on the non-target concept). Such a negative impact on non-target concept generation intensifies with increasing  $s$ , since a larger  $s$  amplifies the token shift. This results in a larger divergence from the original token direction, and eventually more noticeable changes in the generated images.

The erasure performance is less sensitive to  $p$ , but it still has some mild impact. For instance, when using a higher value of  $s$ , a lower  $p$  can mitigate changes in the generated visual content that is relevant to the non-target concepts. This is demonstrated in the bottom-right part of Fig. 8. When  $\epsilon$  decreases to 0.7, setting  $p$  to 40 results in less deviation from the original images as compared to other values. On the other hand, when  $s = 1$ , a higher  $p$  positively affects the preservation of some non-target concepts that are related. This is shown in the bottom-left part of Fig. 8. When  $\epsilon = 0.6$ , the deviation from the original image decreases as  $p$  increases from 40 to 100.

## D. Additional Single-concept Experiments

### D.1. Extended Quantitative Results on Instance and Art Style Erasure

We present the extended quantitative results on instance erasure and art style erasure in Table 4 and 5, respectively. In addition to the CS for the target concept and FID for non-target concepts, we also include the FID for the target concept and CS for non-target concepts. Specifically, FID measures the distribution distance of generated images aligned with the target concept before and after concept erasure, while CS evaluates the semantic consistency between the text prompt of the non-target concept and the generated image after erasure.

However, a lower FID for the target concept only indicates significant visual changes in the generated images but does not confirm that the semantics aligned with the target concept have been fully eliminated. Similarly, a higher CS for the



Figure 8. **Impact of hyperparameter settings on erasure efficacy and prior preservation.** To evaluate how hyperparameters affect this balance, we visualize images generated by AdaVD under various hyperparameter settings for the target concept “*Snoopy*” and the related but non-target concept “*Mickey*”.

non-target concept suggests that the generated image after concept erasure still aligns closely with the text prompt, but does not guarantee small pixel-level changes. In summary, FID for the target concept and CS for the non-target concept cannot directly measure the effectiveness of erasure or prior preservation. Nevertheless, they remain valuable for further verifying and comparing the erasure efficacy and prior preservation.

## D.2. On Celebrity Erasure

We experiment with erasing different celebrity concepts, including “*Bruce Lee*”, “*Marilyn Monroe*”, and “*Melania Trump*”. Five types of prompts were tested, each containing a distinct concept from “*Bruce Lee*”, “*Marilyn Monroe*”, “*Melania Trump*”, “*Anne Hathaway*” and “*Tom Cruise*”. As reported in Table 6, AdaVD consistently exhibits superior erasing efficacy with prior preservation. When erasing different celebrities, AdaVD achieves the lowest or near-lowest CS and FID values, particularly excelling in FID. Although SPM ranks the second in prior preservation based on its FID scores, it falls significantly behind in its overall prior preservation quality, as compared to AdaVD.

Fig. 9 illustrates and compares generated images of methods, where consistent superior performance of AdaVD can be observed. For the target concept “*Marilyn Monroe*”, AdaVD, SPM, and MACE can all successfully remove the celebrity identity. But SPM is overly aggressive at erasing, obscuring the facial outlines. For non-target concepts, all the four competing methods have caused some quite strong deviations, altering the original images. This is particularly noticeable in the

	Snoopy		Mickey		Spongebob		Pikachu		Dog		Legislator	
	CS	FID	CS	FID	CS	FID	CS	FID	CS	FID	CS	FID
SD v1.4	28.49	-	26.50	-	27.30	-	27.41	-	24.27	-	23.73	-
Erase <i>Snoopy</i>												
	CS ↓	FID	CS	FID ↓	CS	FID ↓	CS	FID ↓	CS	FID ↓	CS	FID ↓
ConAbl	25.38	103.80	26.68	38.44	27.02	41.59	27.57	29.68	24.12	27.76	23.48	27.36
MACE	20.78	169.22	22.95	118.01	23.33	111.90	25.77	81.99	23.96	43.27	22.25	65.97
SPM	23.89	122.63	26.66	<u>33.06</u>	27.12	<u>34.70</u>	27.51	<u>23.89</u>	24.24	<u>19.61</u>	23.70	<u>18.26</u>
NP	23.66	125.98	26.14	59.58	26.66	78.74	27.36	52.37	23.89	67.51	22.16	55.22
SLD	27.84	64.78	26.46	48.12	27.52	55.36	27.33	38.74	24.03	41.95	22.80	49.08
Ours	<b>20.28</b>	120.46	26.53	<b>5.72</b>	27.25	<b>8.56</b>	27.40	<b>5.79</b>	24.27	<b>2.32</b>	23.77	<b>6.07</b>
Erase <i>Snoopy and Mickey</i>												
	CS ↓	FID	CS ↓	FID	CS ↓	FID ↓	CS	FID ↓	CS	FID ↓	CS	FID ↓
ConAbl	24.26	119.96	24.08	96.94	27.02	46.32	27.75	39.63	23.98	30.57	23.33	27.49
MACE	20.74	171.16	<u>20.71</u>	140.50	25.87	51.49	25.87	110.67	23.82	52.07	21.70	77.13
SPM	23.16	128.08	22.81	115.02	26.92	<u>41.58</u>	27.45	<u>31.77</u>	24.13	<u>21.96</u>	23.60	<u>23.69</u>
NP	23.59	124.10	24.85	83.68	26.69	81.41	27.27	50.10	23.62	65.93	21.84	58.88
SLD	27.76	59.97	26.74	50.16	27.53	54.59	27.29	39.24	23.97	41.62	22.66	50.13
Ours	<b>20.29</b>	121.12	<b>19.93</b>	108.22	27.27	<b>9.34</b>	27.42	<b>5.84</b>	24.26	<b>2.41</b>	23.73	<b>6.43</b>
Erase <i>Snoopy and Mickey and Spongebob</i>												
	CS ↓	FID	CS ↓	FID	CS ↓	FID ↓	CS	FID ↓	CS	FID ↓	CS	FID ↓
ConAbl	23.94	126.70	23.64	105.07	25.04	108.67	27.76	51.20	23.83	23.83	23.17	30.03
MACE	20.48	172.80	<u>20.50</u>	143.66	21.59	120.87	24.38	99.68	23.70	47.46	21.74	70.38
SPM	22.81	133.06	22.35	121.85	<u>20.82</u>	152.72	27.45	39.83	24.10	<u>22.68</u>	23.52	<u>25.31</u>
NP	24.29	129.75	24.76	89.74	25.31	106.30	27.28	64.75	23.55	65.10	21.63	59.33
SLD	27.84	58.16	26.71	49.70	27.60	54.61	27.35	39.41	23.90	42.32	22.46	49.88
Ours	<b>19.39</b>	124.49	<b>19.73</b>	112.97	<b>20.34</b>	118.47	27.42	<b>6.85</b>	24.27	<b>2.79</b>	23.76	<b>7.26</b>

Table 4. **Extended quantitative comparison of single- and multi-instance erasure.** The best and second-best results are marked in **bold** and underlined, respectively. Columns in gray indicate items that do not directly reflect erasure efficacy or prior preservation performance.



Figure 9. **Qualitative comparison of celebrity erasure.** Our AdaVD can effectively remove the target concept “Marilyn Monroe” while preserving non-target celebrities like “Bruce Lee” and “Melania Trump”.

generated images from the prompt corresponding to “Melania Trump”. For instance, MACE and SPM have introduced an additional arm in the left image, NP has altered the original pose, and SLD has caused a severe visual change in the mouth and eye areas. In contrast, AdaVD is able to successfully maintain all the non-target images with minimal visual changes.



	Van Gogh		Picasso		Monet		Andy Warhol		Caravaggio	
	CS	FID	CS	FID	CS	FID	CS	FID	CS	FID
SD v1.4	29.20	-	28.84	-	29.41	-	29.73	-	27.09	-
<i>Erase Van Gogh</i>										
	CS ↓	FID	CS	FID ↓	CS	FID ↓	CS	FID ↓	CS	FID ↓
ConAbl	28.80	120.93	28.10	71.71	25.99	138.72	29.34	70.30	26.83	73.10
MACE	27.74	144.75	28.37	65.77	29.48	69.79	29.30	83.37	27.11	75.41
SPM	<b>24.78</b>	185.50	28.34	<u>62.25</u>	29.34	<u>32.27</u>	29.52	<u>58.30</u>	27.01	<u>61.50</u>
NP	24.90	193.24	25.11	141.56	26.08	124.52	27.06	127.85	25.34	136.32
SLD	27.48	133.07	26.89	103.96	27.61	109.11	28.24	103.89	25.82	119.32
SAFREE	25.82	183.06	25.84	130.35	27.15	128.71	27.20	127.72	25.53	134.46
Ours	<u>24.87</u>	188.94	28.80	<b>6.82</b>	29.43	<b>2.66</b>	29.74	<b>8.36</b>	27.09	<b>6.84</b>
<i>Erase Picasso</i>										
	CS	FID ↓	CS ↓	FID	CS	FID ↓	CS	FID ↓	CS	FID ↓
ConAbl	29.46	58.62	27.72	121.45	26.37	140.34	29.51	73.35	27.17	67.44
MACE	29.73	60.46	27.11	131.82	29.44	49.92	29.65	76.10	27.08	72.85
SPM	29.26	38.79	26.69	157.32	29.44	<u>7.76</u>	29.67	52.00	27.08	51.40
NP	29.28	111.35	<b>26.14</b>	169.23	29.34	91.11	28.14	116.24	26.50	121.82
SLD	29.36	98.21	27.03	105.37	29.79	93.01	28.80	97.00	26.42	110.05
SAFREE	29.96	117.32	26.42	183.80	29.45	93.51	27.88	122.89	26.32	116.51
Ours	29.17	<b>5.49</b>	26.99	132.64	29.43	<b>2.33</b>	29.72	<b>9.38</b>	27.09	<b>7.05</b>
<i>Erase Monet</i>										
	CS	FID ↓	CS	FID ↓	CS ↓	FID	CS	FID ↓	CS	FID ↓
ConAbl	25.84	141.52	25.47	132.10	<u>24.53</u>	143.48	26.25	208.38	25.48	186.26
MACE	29.47	76.90	28.56	69.35	26.89	109.58	29.34	88.35	26.75	81.72
SPM	29.19	<u>41.03</u>	28.65	<u>29.71</u>	27.00	105.09	29.65	<u>31.90</u>	29.65	<u>25.99</u>
NP	26.31	137.21	25.59	126.75	<b>24.47</b>	140.92	27.05	127.22	24.85	135.83
SLD	28.22	94.48	27.10	92.88	25.73	120.14	28.34	100.90	25.45	114.87
SAFREE	26.07	125.98	26.25	119.19	25.33	153.96	26.82	125.27	25.45	129.07
Ours	29.19	<b>6.94</b>	28.80	<b>6.50</b>	26.30	114.06	29.76	<b>8.46</b>	27.10	<b>7.19</b>

Table 5. **Extended quantitative comparison of art style erasure.** AdaVD achieves a superior balance between erasure efficacy and prior preservation, especially excelling in prior preservation. Notably, it outperforms the concurrent method SAFREE, which also employs orthogonal decomposition.

### D.3. On NSFW Erasure

Unlike the erasure of specific instances, art styles, and celebrities, NSFW concept erasure is more challenging. One reason is that the NSFW concepts are often implicit and hidden within prompts that can be particularly rich in their semantics. Also, many NSFW concepts have synonyms, and it is important to remove both the target concept and its synonyms. For instance, when targeting at removing the “*nudity*” concept, it is essential to also remove the “*sexual*” concept. We experiment with erasing the “*nudity*” concept using the I2P benchmark. To examine how well the “*nudity*” concept is erased, we employ the NudeNet with a threshold of 0.3 to detect nudity in the generated images and analyze the total number of nude items and the overall nude images that are detected.

Results are reported in Fig. 10, where, despite the challenges, AdaVD demonstrates a superior nudity erasure performance, with a semi-threshold and a slower increasing rate. It outperforms both training-based and training-free methods, achieving the best or close-to-best success rate in nearly all categories, with approximately 85% of the nude items successfully removed. It is worth mentioning that NudeNet can be overly aggressive at detecting nude items, resulting in detection errors. For example, it may incorrectly classify a circle with a dot as “*Female Breast Exposure*” or a person opening their mouth as “*Male Genitalia Exposure*”. We increased the NudeNet threshold to 0.3, in order to mitigate this issue, but still, there is a detection error. Being examined by an overly strict nudity detector that can flag sometimes healthy or irrelevant content as nude ones, AdaVD achieves the highest erasure rate for nearly all tested nude items compared to other competing methods, as shown in Fig. 10.

	Bruce Lee		Marilyn Monroe		Melania Trump		Anne Hathaway		Tom Cruise	
	CS	FID	CS	FID	CS	FID	CS	FID	CS	FID
SD v1.4	30.77	-	27.70	-	29.80	-	31.96	-	31.12	-
<i>Erase Bruce Lee</i>										
	CS ↓	FID	CS	FID ↓	CS	FID ↓	CS	FID ↓	CS	FID ↓
ConAbl	31.35	87.57	28.23	57.79	29.77	40.95	29.77	40.95	30.97	53.53
MACE	25.04	131.29	28.13	74.80	30.07	68.83	31.91	75.05	28.13	71.20
SPM	27.75	123.67	27.71	26.89	29.81	7.83	31.96	9.46	31.13	28.54
NP	24.70	150.85	26.84	102.67	28.94	82.13	30.34	89.60	29.67	89.92
SLD	28.22	102.26	26.29	87.15	29.43	84.32	30.97	85.37	29.32	94.07
Ours	<b>20.67</b>	138.70	27.70	<b>6.68</b>	29.82	<b>5.08</b>	31.97	<b>6.39</b>	31.10	<b>13.11</b>
<i>Erase Marilyn Monroe</i>										
	CS	FID ↓	CS ↓	FID	CS	FID ↓	CS	FID ↓	CS	FID ↓
ConAbl	30.88	66.97	28.75	88.45	29.69	51.52	32.05	58.57	31.10	54.13
MACE	31.30	76.23	<b>19.52</b>	148.34	31.93	71.05	30.16	74.90	31.52	73.06
SPM	30.76	<u>32.70</u>	21.87	145.81	29.83	<u>25.27</u>	31.96	<u>22.86</u>	31.10	<u>19.34</u>
NP	29.50	113.12	25.86	149.95	29.29	87.27	29.42	98.86	30.02	86.70
SLD	29.59	87.83	26.70	98.51	28.81	107.42	29.25	102.13	30.35	81.12
Ours	30.73	<b>7.88</b>	<u>19.87</u>	116.94	29.80	<b>4.46</b>	31.93	<b>5.43</b>	31.13	<b>9.33</b>
<i>Erase Melania Trump</i>										
	CS	FID ↓	CS	FID ↓	CS ↓	FID	CS	FID ↓	CS	FID ↓
ConAbl	30.62	54.46	28.14	59.10	29.89	79.04	31.94	58.65	31.00	54.50
MACE	31.30	78.07	27.84	71.34	20.71	122.42	31.94	73.49	31.41	71.09
SPM	30.79	<u>14.08</u>	27.63	<u>30.40</u>	<b>23.12</b>	129.68	31.86	<u>28.85</u>	31.10	<u>22.35</u>
NP	29.38	115.35	27.63	103.83	23.73	131.73	28.72	106.04	30.27	106.00
SLD	29.55	90.69	26.24	93.93	25.45	103.52	28.43	104.48	30.47	88.31
Ours	30.75	<b>7.32</b>	27.69	<b>6.86</b>	<u>23.28</u>	96.66	31.95	<b>6.52</b>	31.08	<b>5.74</b>

Table 6. **Quantitative comparison of celebrity erasure.** Compared to both training-based and training-free methods, AdaVD achieves an optimal balance between erasure efficacy and prior preservation, demonstrating exceptional performance, particularly in prior preservation.

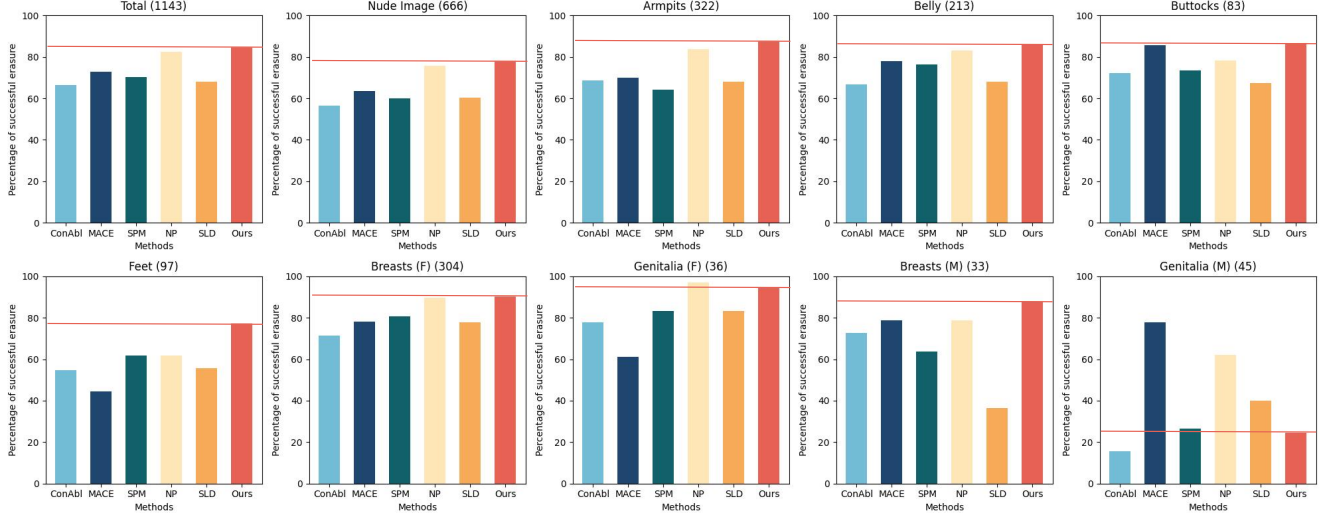


Figure 10. **Performance of AdaVD on NSFW erasure.** The number following each category represents the number of nude items generated by SD v1.4, while each bar illustrates the success rate of erasing the corresponding nude items for each method. Our AdaVD demonstrates superior performance on NSFW erasure, outperforming both training-based and training-free methods.

#### D.4. More Erasure Examples

We demonstrate additional examples for erasing single concepts from prompts that contain such concepts. The experimented concepts include the specific instances of “Statue of Liberty”, “BB8”, “C3PO”, and “Grumpy Cat”, the celebrity “Benicio Del Toro”, and the art style “Cyberpunk”. Among these, “BB8” and “C3PO” are fictional characters, while “Statue of

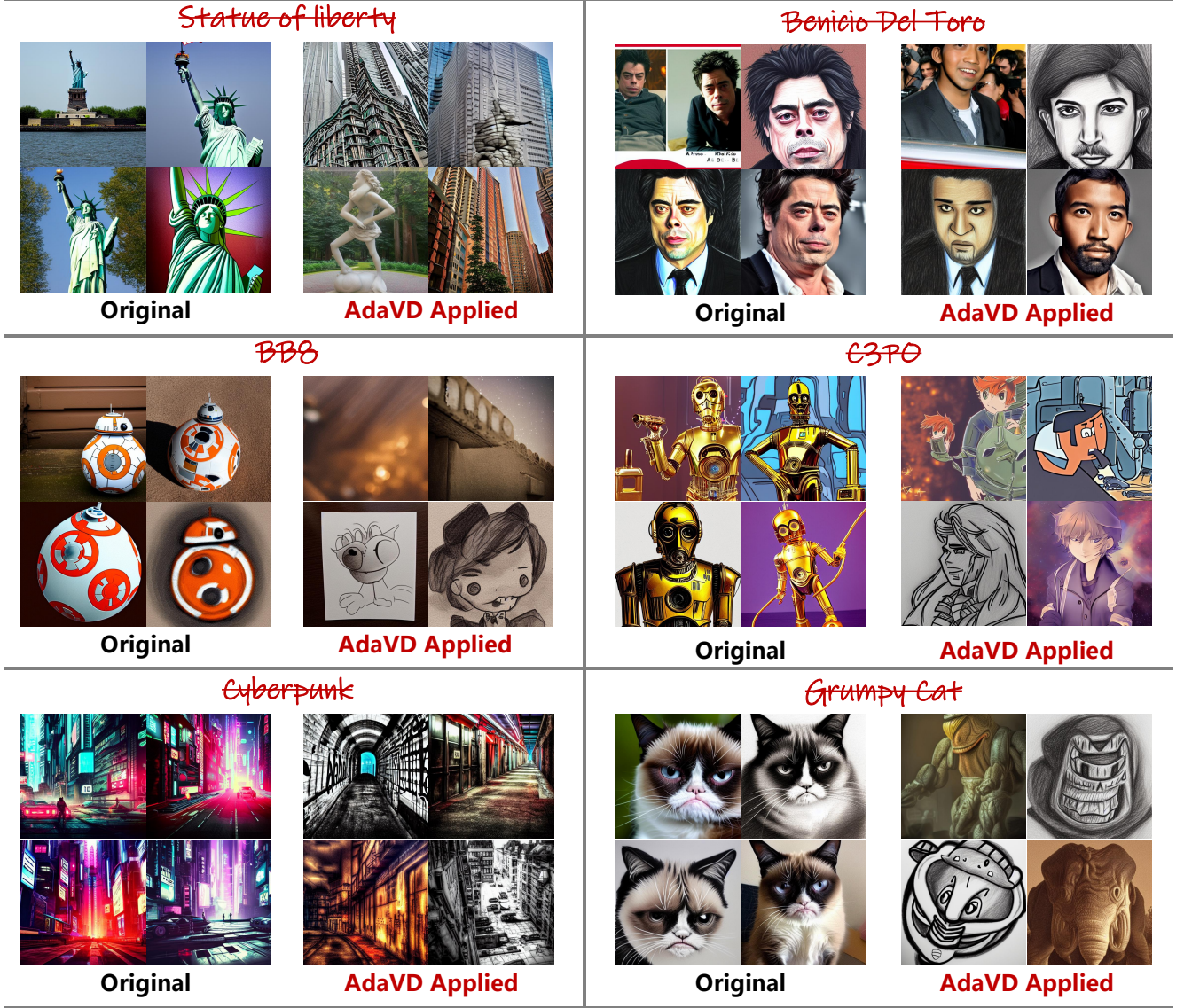


Figure 11. **Extended results of AdaVD in single-concept erasure task.** We present additional generated images after applying AdaVD with SD v1.4 to erase a single concept, further validating the erasure efficacy of our AdaVD.

*Liberty*” and *“Grumpy Cat*” represent realistic entities from daily life. Fig. 11 presents the generated image examples. It can be seen that our AdaVD consistently exhibits superior erasure efficacy across all these concepts, being robust in erasing diverse types of concepts.

## E. On Transferability to Other T2I Models

The proposed AdaVD is a flexible concept erasure approach that can be transferred to other T2I diffusion models. In addition to SD v1.4, as experimented in the main paper, we conduct additional experiments to demonstrate its transferability and effectiveness by integrating it with a series of other T2I diffusion models.

### E.1. AdaVD on SDXL v1.0

We integrate AdaVD with SDXL v1.0 [31] which has a different architecture from SD v1.4-v2.1. It employs two distinct text encoders to process textual prompts, and their outputs are concatenated and fed into the CA layers to interact with the latent representations of the noisy images. Also, the generated text embeddings are enhanced by time embeddings to ensure the



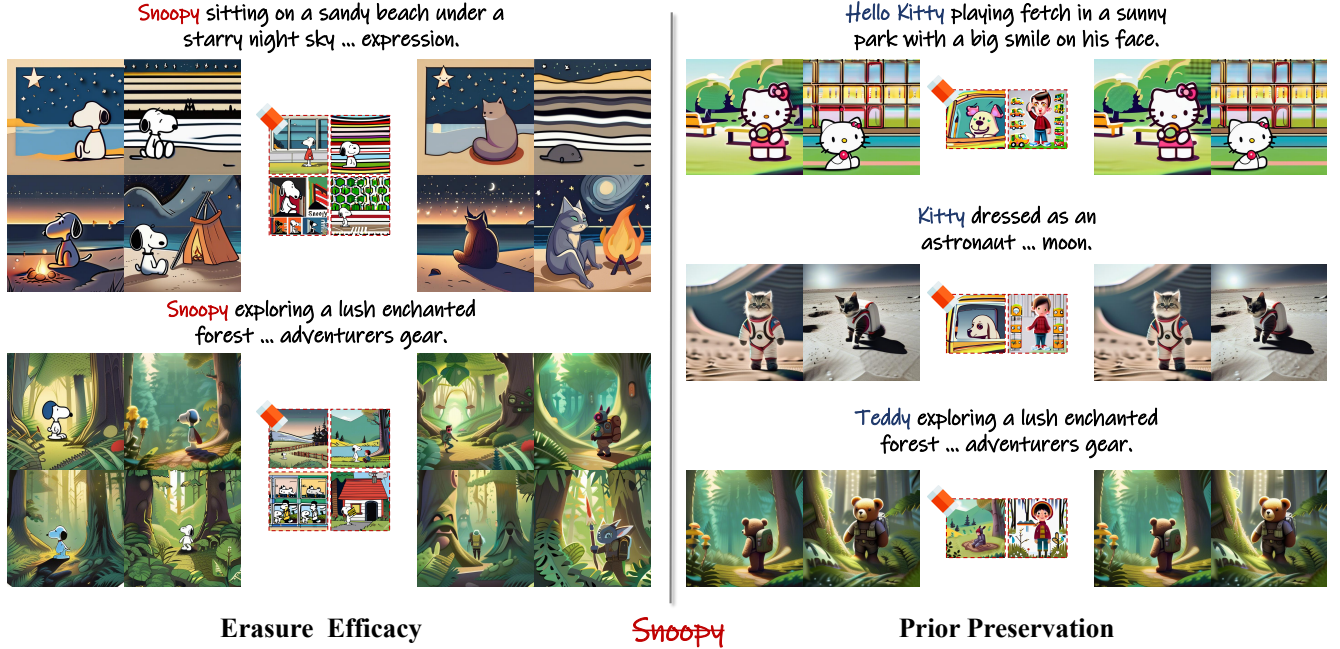


Figure 12. **Results of AdaVD on SDXL v1.0 for erasing “Snoopy”**: Our AdaVD effectively supports SDXL v1.0, which has a different structural design than SD v1.4, in achieving effective erasure of the target concept. Additionally, AdaVD demonstrates excellent prior preservation, as evidenced by its ability to generate non-target concepts like “Hello Kitty”, “Kitty”, and “Teddy” even with semantically rich prompts. AdaVD successfully retains nearly all details in non-target content, underscoring its capability for precise erasure without compromising unrelated elements.

alignment between textual prompts and timesteps. Following the same approach as how it is coupled with SD v1.4, AdaVD is applied in the value space at each CA layer within the UNet of SDXL v1.0. For the target concepts, both sets of their embeddings computed by the two text encoders are pre-processed following the procedure outlined in Sec. 3.2, then they are used to start the erasure process following the method outlined in Sec. 3.3.1.

Fig. 12 demonstrates the generated image examples by coupling AdaVD with SDXL v1, for long and semantically rich prompts that (do not) contain the “Snoopy” concept while with the target concept “Snoopy” to erase. Although the prompts are more complex, they do not appear challenging for AdaVD to handle. AdaVD can still accurately identify and extract the relevant semantic components associated with the target concept, and can precisely erase these without affecting the background generation. We visualize the erased component for each generated image in the smaller images within each example block of Fig. 12, following the same approach as explained in the 2nd paragraph of Section 4.4. These serve as supporting evidence, showing what semantic content has been removed by AdaVD. For those prompts containing only the non-target concepts, AdaVD successfully retains nearly all the details of the non-target content, producing images that are virtually identical to those generated by the original SDXL v1.0.

## E.2. AdaVD on SDv3

A growing trend in text-to-image generative diffusion models is replacing U-Net with DiT as the noise predictor. Different from U-Net, DiT uses a transformer-based architecture, enhancing scalability in image generation. To validate the performance of our AdaVD in DiT-based diffusion models, we conduct experiments on SDv3. Different from SDv1.4 and SDXL, SDv3 uses the T5 text encoder [34], alongside other encoders, to generate text embeddings for image generation. During the target embedding pre-processing phase, we handle text embeddings differently depending on the encoder: for embeddings from the CLIP text encoder, we replicate the last subject token, while for those from T5, we spread the mean embedding of all real word tokens. As shown in Fig. 13, SDv3 successfully removes the target concept “Snoopy” during the generation process while preserving the integrity of non-target concepts such as “Stitch”, “Mickey”, and “Spongebob”. This highlights the strong prior preservation capability of AdaVD.



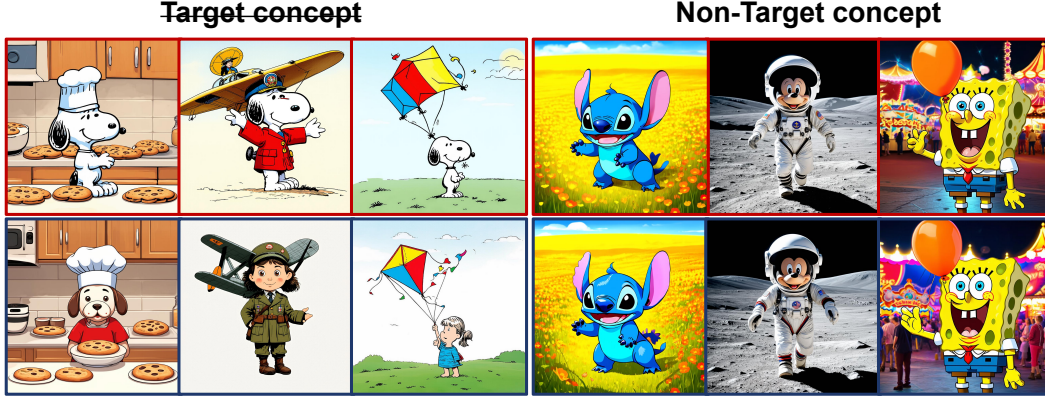


Figure 13. **Results of AdaVD on SDv3 for erasing “Snoopy”**: The images with red and blue borders represent the before and after concept erasure, respectively. Our AdaVD effectively enables SDv3 to erase the target concept “Snoopy” while preserving other semantic elements in the generated images. Moreover, AdaVD demonstrates outstanding prior preservation by ensuring that non-target concepts such as “Stitch”, “Mickey”, and “Spongebob” remain highly similar to the generated images before concept erasure.

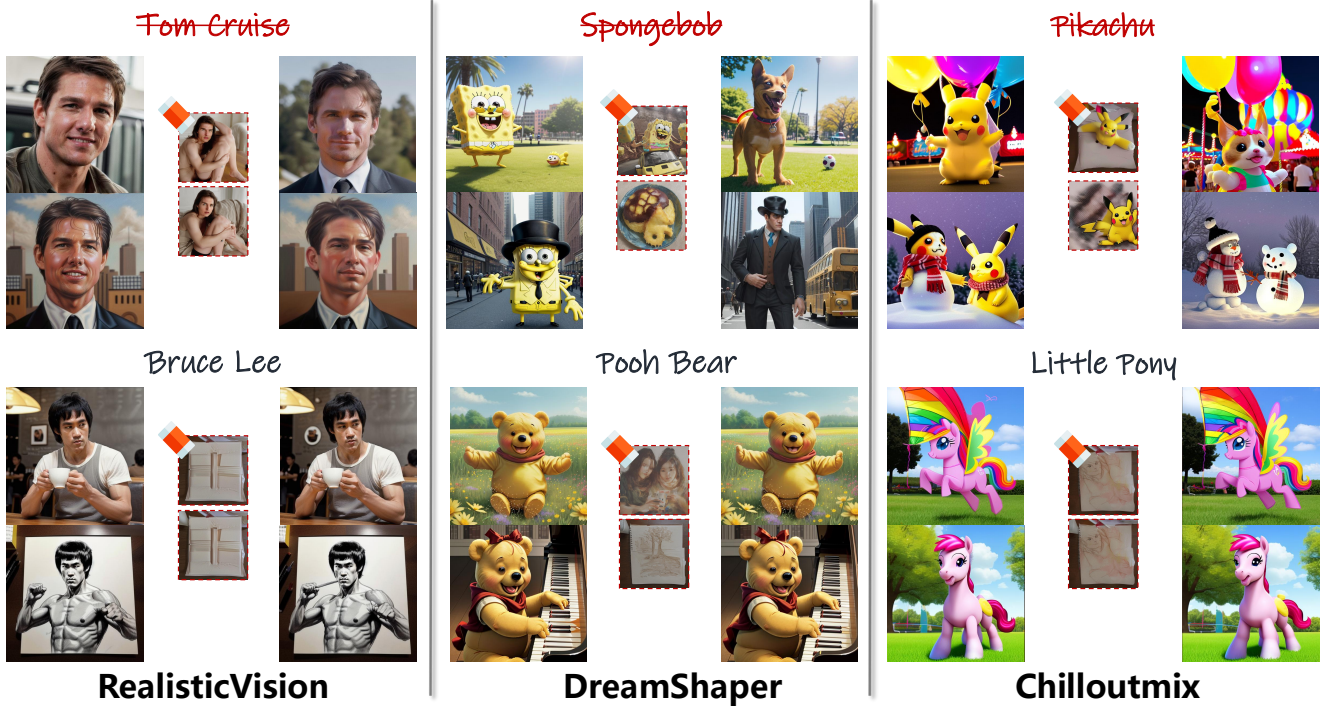


Figure 14. **Results of AdaVD on other SD versions**. Our AdaVD illustrates a high performance of both erasure efficacy and the prior preservation across SD with difference versions and easing different concepts.

### E.3. AdaVD on Community SD Versions

We also couple AdaVD with several community versions of SD, including RealisticVision [7], Dreamshaper [6], and Chilloutmix [5], which are all fine-tuned based on SD v1.5. These versions target high-quality image generation with specific generation objectives. For example, RealisticVision specializes in generating lifelike images, while Dreamshaper excels in producing highly imaginative visuals. We experiment with removing the target concept “Tom Cruise” from the text prompt corresponds to “Tom Cruise” and “Bruce Lee” for RealisticVision, removing “Spongebob” from the text prompt corresponds to “Spongebob” and “Pooh Bear” for Dreamshaper, and removing “Pikachu” from the text prompt corresponds to “Pikachu” and “Little Pony” for Chilloutmix.

Fig. 14 presents the generated image examples. The results show that AdaVD is capable of effectively erasing the target

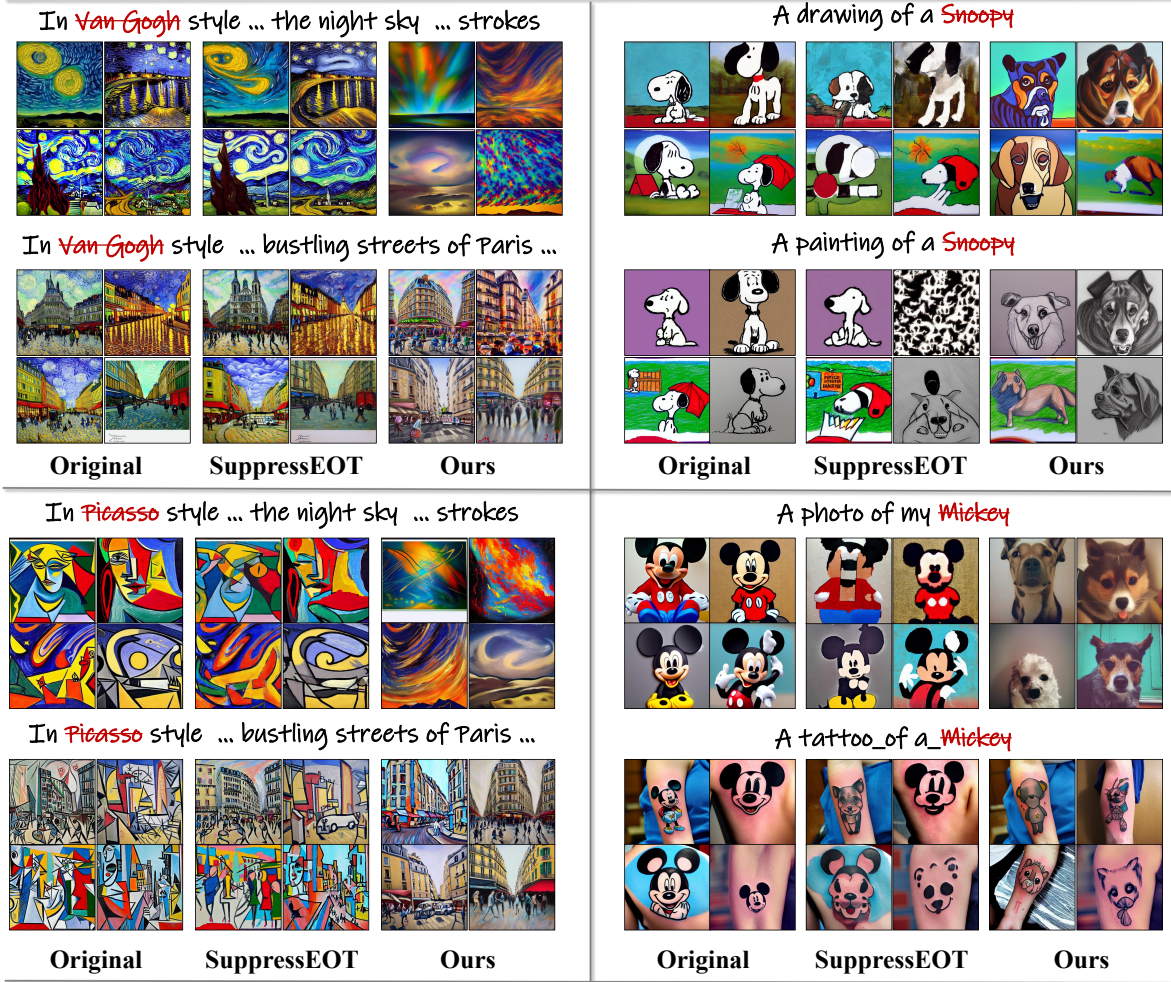


Figure 15. **Qualitative comparison between SuppressEOT and AdaVD.** We compare our AdaVD with SuppressEOT in single instance concept and art style erasure, demonstrating that AdaVD achieves more precise and effective erasure.

concept while preserving the integrity of the non-target content. For all the experimented community versions, AdaVD can precisely locate the semantic space aligned with the target concept and isolate it with minimal disruption to the non-target semantics. Fig. 14 also visualizes the erased components as the smaller images within each example block, as in Fig. 12. Overall, the visualized erased components for prompts containing the target concepts show a high similarity to the target semantics. In contrast, for prompts corresponding to non-target concepts, the erased components lack meaningful semantic information. These serve as additional evidence, showing the effectiveness of AdaVD.

## F. Comparison with Additional Baselines

### F.1. Comparison with SAFREE

Orthogonal complement is widely used to decouple and separate out unwanted information. The art of using the orthogonal complement for concept erasure is on designing/deciding what space/direction to apply orthogonal complement, how to adjust removal strength, how to embed orthogonal complement in an algorithm to optimize its effect, etc. There is a concurrent work, SAFREE [51], which also used orthogonal complement to facilitate concept erasure, but in completely different ways. Our AdaVD performs orthogonal complement in value spaces of attention layers within a diffusion model. Due to its effectiveness, there is no need for any complementary design, but a soft control of removal strength through a shift factor. Different from our AdaVD, SAFREE performs orthogonal complement over diffusion model input, i.e., text embedding space. This approach necessitates complementary design elements, such as masking, another projection, and modifying



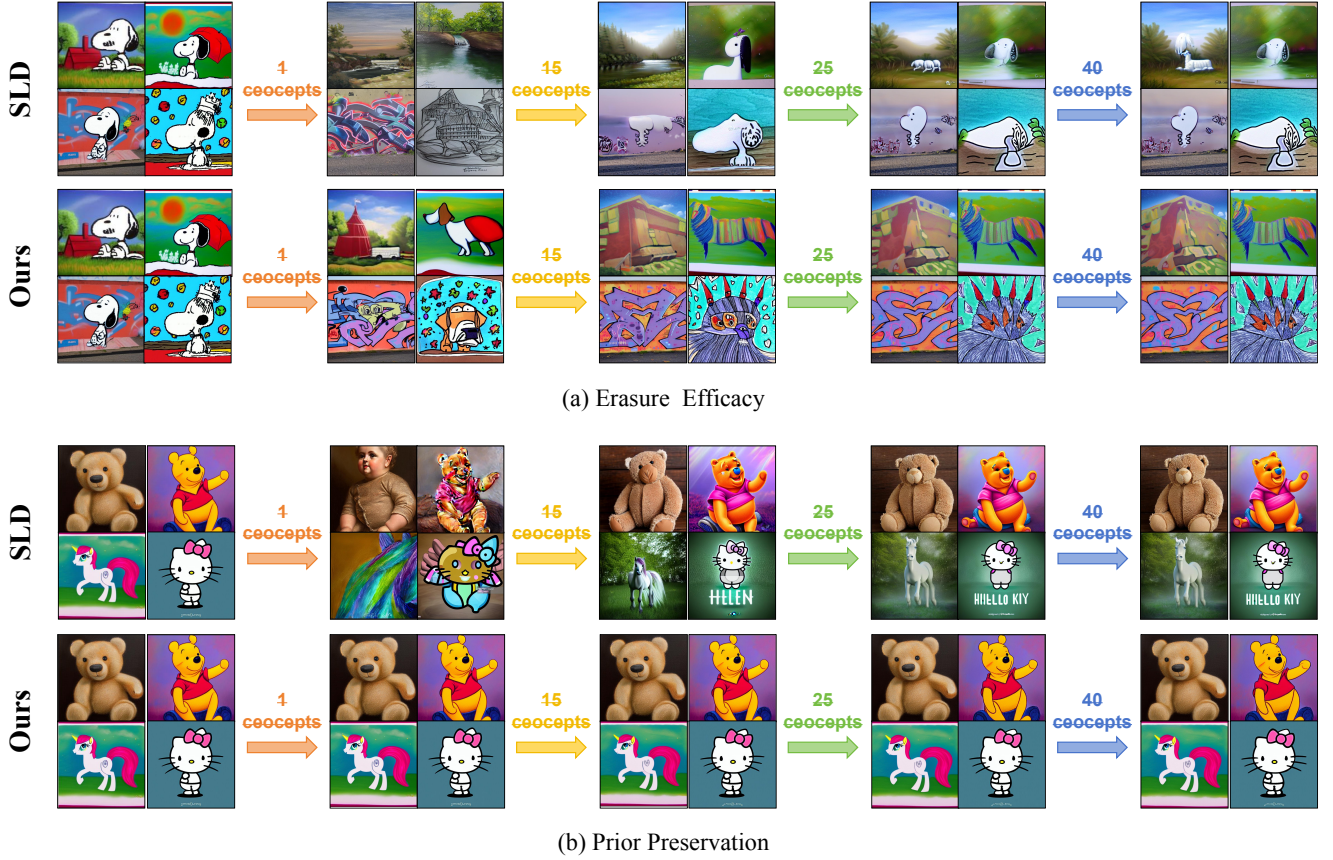


Figure 16. **Examples of generated images for multi-concept erasure.** The illustrated examples show a consistently high performance of AdaVD in both erasure efficacy and prior preservation as the number of erased concepts increases, as compared to SLD.

detoxified embedding by Fourier transform. To control removal strength, it also uses a hard selection of whether to adopt the final detoxified embeddings. We also conduct experiments to compare the performance of AdaVD and SAFREE in erasing art style concepts. As shown in Table 5, AdaVD achieves excellent prior preservation performance and second-best erasure efficacy, outperforming SAFREE, especially in prior preservation.

## F.2. Comparison with SuppressEOT

In this additional experiment, we compare with a special concept erasure method SuppressEOT [24], which requires the users to specify the positions of the erased concepts within the prompt. Because of this user-involved setting, SuppressEOT is only applicable to specific prompts and is unable to achieve system-wide concept erasure. Therefore, we only conduct a qualitative comparison of the erasure efficacy. Results are reported in Fig. 15, where a comparison of art style erasure is shown on the left side, while the instance erasure results are displayed on the right.

It can be seen from Fig. 15 that AdaVD is precise and effective in erasing various concepts, achieving significantly higher erasure performance across a diverse range of use cases. Unfortunately, SuppressEOT consistently fails to remove completely the target concept from the generated images. For instance, when erasing “Mickey” and “Snoopy”, SuppressEOT is not even able to erase the general outlines of these specific instances. The unsatisfactory erasure performance of SuppressEOT likely stems from the fact that it was originally designed for image editing rather than concept erasure. In image editing scenarios, preserving all the details of a prompt except for the target concept is important. This is different from the requirement of concept erasure, where the prior preservation is needed only for the generation of non-target content. Such a difference in design requirement can inherently compromise the erasure efficacy of SuppressEOT.

Number	Target Concepts
1	<i>Snoopy</i>
15	<i>Snoopy, Mickey, Crystal, Pikachu, Legislator, Bruce Lee, Marilyn Monroe, Tom Cruise, Anne Hathaway, Melania Trump, Van Gogh, Picasso, Rembrandt, Andy Warhol, Caravaggio</i>
25	<i>Snoopy, Mickey, Crystal, Pikachu, Legislator, Bruce Lee, Marilyn Monroe, Tom Cruise, Anne Hathaway, Melania Trump, Van Gogh, Picasso, Rembrandt, Andy Warhol, Caravaggio, Samoyed, Doraemon, Tom, Adam Driver, Adriana Lima, Amber Heard, Amy Adams, Andrew Garfield, Angelina Jolie, Anjelica Huston</i>
40	<i>Snoopy, Mickey, Crystal, Pikachu, Legislator, Bruce Lee, Marilyn Monroe, Tom Cruise, Anne Hathaway, Melania Trump, Van Gogh, Picasso, Rembrandt, Andy Warhol, Caravaggio, Samoyed, Doraemon, Tom, Adam Driver, Adriana Lima, Amber Heard, Amy Adams, Andrew Garfield, Angelina Jolie, Anjelica Huston, Bradley Cooper, Bruce Willis, Bryan Cranston, Cameron Diaz, Channing Tatum, Charlie Sheen, Charlize Theron, Chris Evans, Chris Hemsworth, Chris Pine, Barack Obama, Beth Behrs, Bill Clinton, Bob Dylan, Bob Marley</i>

Table 7. **Number of concepts to be erased and their corresponding lists.** The number of concepts ranges from 1 to 40, demonstrating the efficacy of AdaVD in handling multi-concept erasure.

## G. Additional Experiments and Analysis on Multi-Concept Erasure

### G.1. On Erasing More Multi-concepts

We conduct additional experiments, investigating how our approach performs as the number of erased concepts increases, under a progressive setting. We evaluate our AdaVD by first erasing one concept “*Snoopy*” and gradually increasing the number of erased concepts to 15, 25, and 40. The details of the concepts to be erased for each case are listed in Table 7. We work with the base T2I model SD v1.4 and compare it with the existing approach SLD. The results are presented in Fig. 16, which extends Fig. 1.

It can be observed from the top erasure efficacy block of Fig. 16 that SLD gradually loses its precision when removing the target concepts. This is possible because SLD concatenates the target concepts into a prompt for guiding the generation process. When erasing too many concepts, the text encoder struggles to focus on each individual concept, resulting in diminished erasure efficacy. Additionally, some concepts may be truncated due to the token length limitation of the text encoder’s tokenizer. Differently, AdaVD achieves consistently high performance in multi-concept erasure. It constructs a value subspace based on the orthogonal complement of all the target concepts, which ensures that no information regarding any individual concept is lost.

The bottom prior preservation block of Fig. 16 shows that AdaVD is able to generate images nearly identical to the original ones, demonstrating a superior performance in prior preservation. But SLD struggles to preserve prior knowledge, for not only the more challenging case of removing a high number of concepts but also the simple case of removing one single concept. It is worth noting that some slight change can be accumulated and amplified as the number of erased concepts increases, as shown in the hands and mouth of the generated image of “*Pooh Bear*” by our AdaVD. Also, small pixel-level changes may grow into catastrophic forgetting with an increasing number of erased concepts due to error accumulation. Therefore it is important to use FID to evaluate the performance of prior preservation, as images that closely match the originals at pixel level should result in a low FID score.

### G.2. On Transferability to Other T2I Models

In this additional experiment, we integrate AdaVD with two other T2I diffusion models, including DreamShaper [6] and RealisticVision [7], assessing its multi-concept erasure performance. Two multi-concept erasure scenarios are experimented with: one is cross-application erasure as described in SPM [27], and the other is multi-instance erasure. Results of the cross-application erasure are presented in the top half of Fig. 17, demonstrating the generated images after erasing “*Snoopy*”, “*Van Gogh*”, and the two concepts together. Results of the multi-instance erasure are shown at the bottom of Fig. 17, demonstrating the generated images after erasing “*Mouse*”, “*Dog*”, and both concepts. Overall, AdaVD achieves a high erasure precision. It can be seen from Fig. 17 that, when aiming at a single concept erasure, other concepts specified in the prompt remain faithfully in the generated image; and when aiming at erasing multiple concepts, all the relevant visual content is also removed successfully. This serves as evidence that AdaVD is capable of a robust and precise erasure.



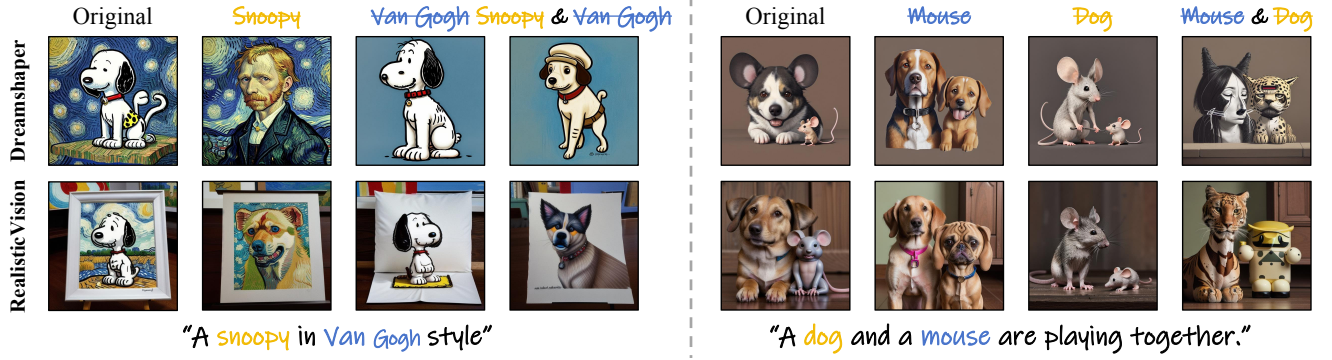


Figure 17. **Results of AdaVD on multi-concept erasure across different SD versions.** We assess the performance of AdaVD on multi-concept erasure across various community versions of SD under diverse erasure scenarios, including cross-application erasure as outlined in SPM [27] and multi-instance erasure. These evaluations further highlight the robustness and effectiveness of AdaVD in addressing the challenges of the multi-concept erasure task.

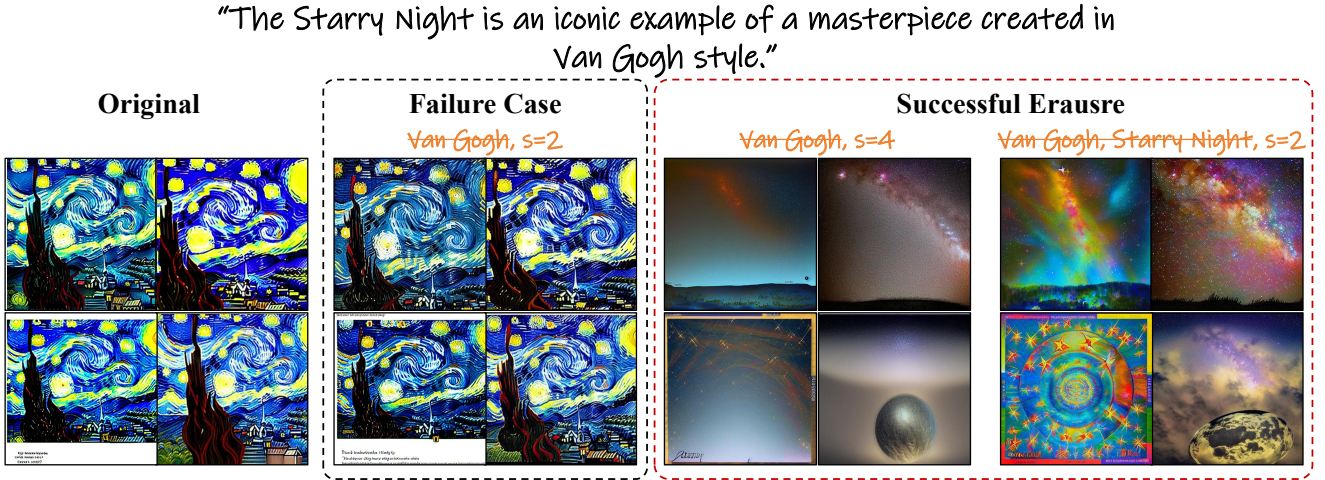


Figure 18. **Failure case** when erasing “Van Gogh” and its solution.

## H. Failure Case Study

Despite its success, there exist concepts that AdaVD struggles to erase. We present a few failure cases in Fig. 18. For instance, it is challenging for AdaVD to erase “Van Gogh” from a prompt like “The Starry Night is an iconic example of a masterpiece created in Van Gogh style.” The challenge is likely to stem from the presence of multiple tokens, *e.g.*, “Starry Night”, that is highly coupled with the target concept. In this case, a small value of the scaling hyper-parameter  $s$  as used by the shift factor in Eq. (6) is insufficient to eliminate effectively the target semantics across all the relevant tokens. Nevertheless, this issue can be mitigated by doubling  $s$  or incorporating additional related target concepts to erase, *e.g.*, “Starry Night”, evidenced by the right side of Fig. 18.