

Unlearning Concepts from Text-to-Video Diffusion Models

Shiqi Liu¹, Yihua Tan^{1*},

¹Huazhong University of Science and Technology, Wuhan 430074, PR China
shiqi.liu647@foxmail.com, yhtan@hust.edu.cn

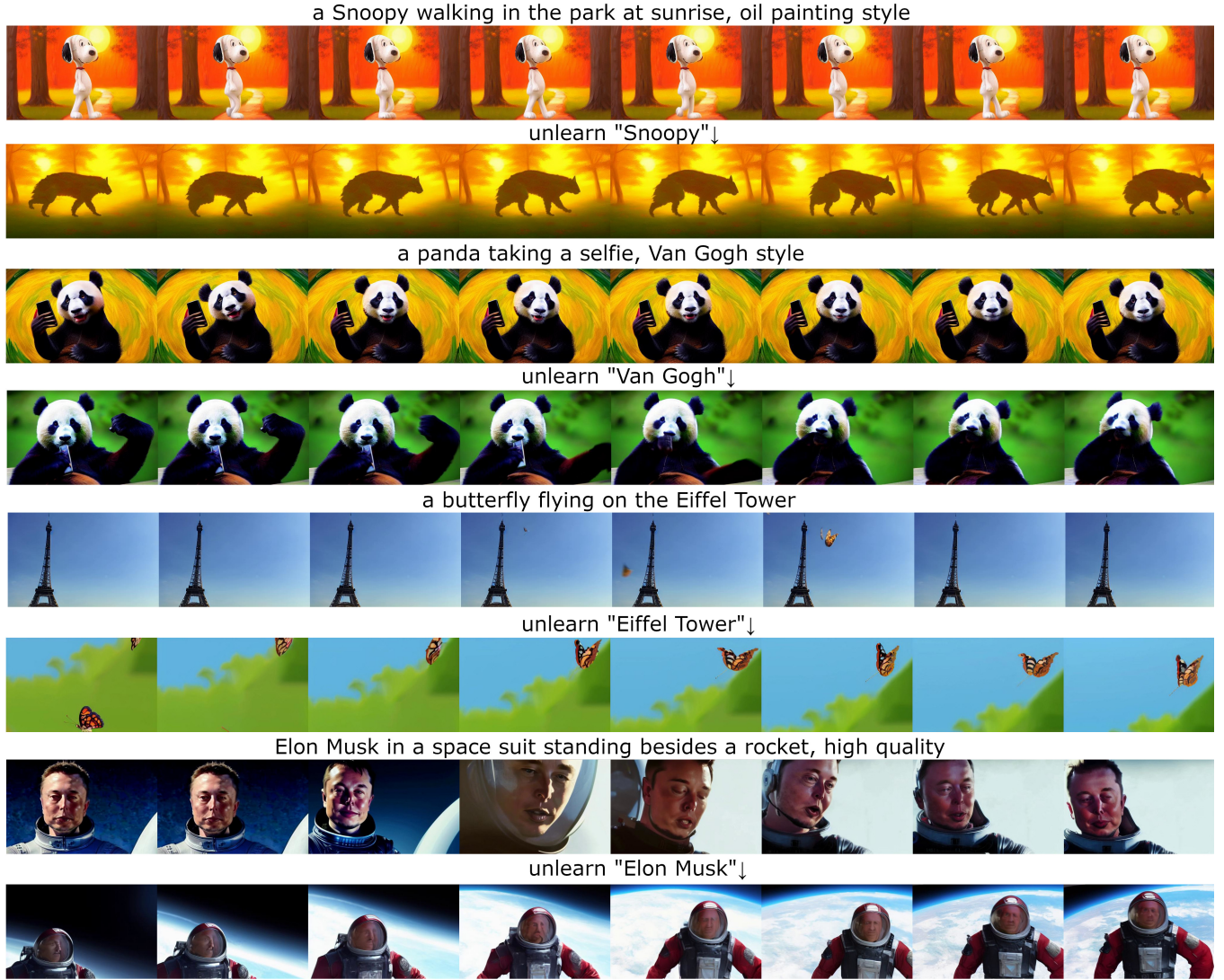


Figure 1: The comparison of the concept-preserved and concept-unlearned videos generated by our algorithm.

Abstract

*Corresponding author.
Copyright © 2025, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

With the advancement of computer vision and natural language processing, text-to-video generation, enabled by text-to-video diffusion models, has become more prevalent. These models are trained using a large amount of data from the internet. However, the training data often contain copyrighted

content, including cartoon character icons and artist styles, private portraits, and unsafe videos. Since filtering the data and retraining the model is challenging, methods for unlearning specific concepts from text-to-video diffusion models have been investigated. However, due to the high computational complexity and relative large optimization scale, there is little work on unlearning methods for text-to-video diffusion models. We propose a novel concept-unlearning method by transferring the unlearning capability of the text encoder of text-to-image diffusion models to text-to-video diffusion models. Specifically, the method optimizes the text encoder using few-shot unlearning, where several generated images are used. We then use the optimized text encoder in text-to-video diffusion models to generate videos. Our method costs low computation resources and has small optimization scale. We discuss the generated videos after unlearning a concept. The experiments demonstrates that our method can unlearn copyrighted cartoon characters, artist styles, objects and people’s facial characteristics. Our method can unlearn a concept within about 100 seconds on an RTX 3070. Since there was no concept unlearning method for text-to-video diffusion models before, we make concept unlearning feasible and more accessible in the text-to-video domain.

Introduction

Recent text-to-video diffusion generative models(Wang et al. 2023; Ho et al. 2022; Yin et al. 2023) have attracted attention because of their outstanding video quality, stable learning procedure, and seemingly infinite generation capabilities, surpassing the previous state-of-art generative adversarial networks(Goodfellow et al. 2020, 2014). Classifier-free guidance(Ho and Salimans 2021) allows us generate high-quality videos on the basis of natural language input. These models are able to imitate a wide range of concepts since they are trained on vast internet datasets.

Their ability to imitate potentially copyrighted content is a major concern regarding text-to-video models. They can faithfully generate copyrighted videos such as “Snoopy,” an iconic beagle dog from the beloved comic strip Peanuts, as shown in Figure 1. The AI-generated art is on par with human-generated art. Another issue is that the models can faithfully replicate an artist’s style. Users of large-scale text-to-video generation systems can use prompts including “in the style of [artist]” to mimic the styles of specific artists, which may reduce the value of the original work. The Van Gogh-style video “a panda taking a selfie” is shown in Figure 1. Some artists have sued the makers and providers of certain generation models, raising new legal issues (Setty 2023).

Apart from copyright infringement issues, privacy and safety are other major concerns. Text-to-video diffusion models can generate specific facial characteristics through text prompts that include names, if the training datasets contain corresponding videos of those names. An example of this is shown in Figure 1: a generated video of Elon Musk. However, this generation of facial characteristics may raise concerns about privacy and portrait rights as outlined in the Civil Code of the People’s Republic of China. Additionally, malicious use of these generated videos could contribute to the spread of fake news and misinformation. Text-to-video

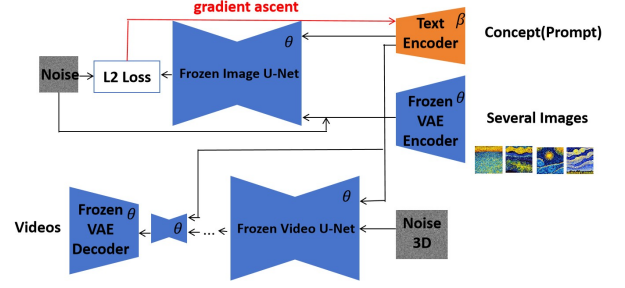


Figure 2: Overview of our proposed method. θ denotes that parameters are fixed and β denotes that parameters need to be optimized.

diffusion models are also capable of generating nude and pornographic videos. These concerns all necessitate technical solutions.

Cleaning the datasets and retraining the text-to-video models require a great amount of work and expenditure. For example, training text-to-video diffusion models, as described in (Wang et al. 2023; Ho et al. 2022), requires about 10 million videos. One feasible approach is to use unlearning methods (Bourtoule et al. 2021), which are proposed to eliminate the influence of specific data or concepts. Previous studies have successfully unlearned specific concepts from text-to-image models by optimizing the weights of U-Net (Ronneberger, Fischer, and Brox 2015), the generative module (Gandikota et al. 2023; Kumari et al. 2023; Gandikota et al. 2024; Zhang et al. 2024; Zhao et al. 2024). A problem with this method is that optimizing the parameters of U-Net can lead to a decline in generation quality. Another alternative is to unlearn specific concepts from text-to-image models by optimizing the parameters of the text encoder (Radford et al. 2021; Raffel et al. 2020). This method (Fuchi and Takagi 2024) uses gradient ascent of parameters with regard to the concept to be unlearned on the images related to the concepts. However, due to the high computational complexity and relative large optimization scale, there is little work on unlearning methods for text-to-video diffusion models. Since some text-to-image diffusion models (Rombach et al. 2022) and text-to-video diffusion models (Wang et al. 2023) share the same text encoder, it is natural to question whether we can transfer the unlearning capability of text-to-image diffusion models to text-to-video diffusion models.

We propose a method, based on text-to-image diffusion models (Fuchi and Takagi 2024), for unlearning specific concepts without altering the video-generating module. We aim to achieve unlearning by altering the embedding of text conditioning in the text encoder. The text-to-video domain unlearning is with high computational complexity and requires relative large optimization scale. We address the problem by using the same text encoder to transfer the unlearning capability of text-to-image diffusion models to text-to-video diffusion models. The transfer learning confines the computational complexity into the text-to-image level and reduces the optimization scale. We utilize a few images of the concept to make small changes to remove the concept from the

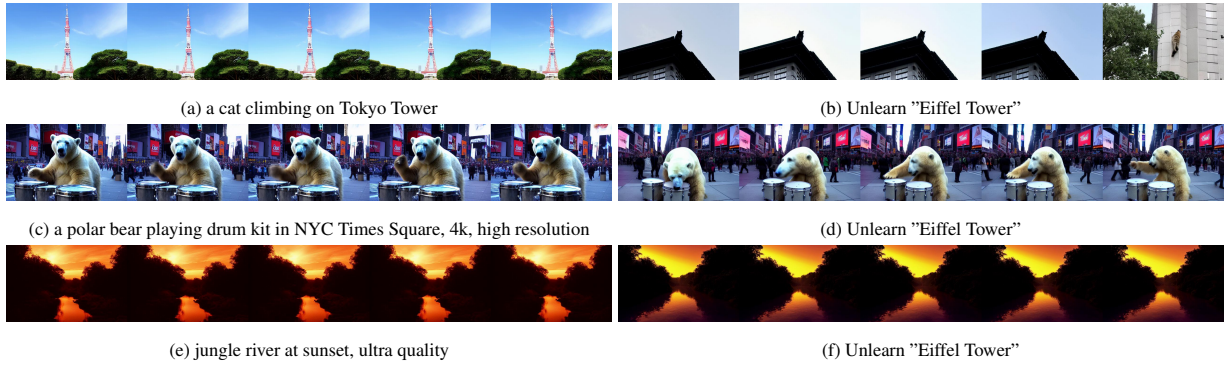


Figure 3: The comparison of the effects on other concepts after unlearning "Eiffel Tower".

text encoder in the image domain and reuse the text encoder in the text-to-video diffusion models to remove the concept in the video domain. Since there are only minor parameter adjustments using image domain optimization on the text encoder, it operates quickly with small costs compared to unlearning in the text-to-video domain. According to Figure 1, it is evident that our method can successfully unlearn concepts.

Our contributions are the following:

- The unlearning procedure transfers from text-to-image diffusion models to text-to-video diffusion models.
- The unlearning procedure takes only about 100 seconds on an RTX 3070.
- Concept unlearning is achieved by providing a few images regarding the concept to be unlearned without using videos or optimizing U-Net.

Related Works

Text-to-Video Diffusion Models and Text Encoder

Due to their simple training objective, the stable training processes, and good generation performance, denoising diffusion models (Ho, Jain, and Abbeel 2020) have been successful in image generation (Rombach et al. 2022) and have gradually attracted more attentions on video generation (Wang et al. 2023). Some representative works include Imagen video (Ho et al. 2022), which uses a cascade sampling pipeline for video generation, NUWA-XL (Yin et al. 2023), which uses a "coarse-to-fine" generation pipeline, and Lavi (Wang et al. 2023), which uses a cascade latent diffusion generation pipeline.

Text-to-video generation utilizes a text encoder to encode the prompt, providing a semantic condition for video generation. There are two commonly used text encoders. One is CLIP (Radford et al. 2021). The other is T5 (Raffel et al. 2020; Ni et al. 2021). CLIP is trained on 400 million (image, text) pairs collected from the internet (Radford et al. 2021). The author compared its performance against over 30 different existing computer vision datasets and the model is competitive with supervised baseline (Radford et al. 2021). T5 performs well on sentence transfer tasks and is state-of-the-art sentence embedding model (Raffel et al. 2020). In many

works, the text encoder is shared and fixed across the video and image domain (Wang et al. 2023) (Rombach et al. 2022).

(Fuchi and Takagi 2024)'s work showed that we can unlearn the concept by adjusting only the parameters of the text encoder in text-to-image diffusion models. Our method builds upon the work of (Fuchi and Takagi 2024). By utilizing the text-to-image and text-to-video models that shared the same text encoder, the unlearning effect on the text encoder on text-to-image domain can be transferred to the text-to-video domain.

Memorization and Unlearning

While the original goal of machine learning is to generalize rather than memorize, large diffusion models are capable of both exact memorization (as shown in (Gu et al. 2023)) and unintentional memorization (Gu et al. 2023). The memorization phenomenon could lead to copyright and privacy issues, prompting the development of unlearning techniques. Unlearning aims to eliminate the influence of specific training samples on the models by adjusting their parameters. This essentially makes the model behave as if it had never encountered those samples. Some unlearning approaches focus on unlearning specific training samples, while others target broader concepts (Ma et al. 2024), as in our work. Our work corresponds to the latter unlearning purpose, aiming to unlearn the high-level concepts such as artist styles, concrete copyright icons, and individuals' appearances.

Unlearning Concepts from Text-to-Image Diffusion and Text-to-Video Diffusion

Many video-generation and image-generation models are trained on a massive amount of data on the Internet (Radford et al. 2021). As a result, these datasets often contain copyrighted material and privacy-sensitive content, which causes problems when deploying the models to the market. Cleaning the data and retraining the model is one method to fix the issue. However, the cost of retraining a model trained on the datasets with millions of samples is often unaffordable. The more practical method is to unlearn some data or concepts of the baseline models (Bourtole et al. 2021).

There are several works on unlearning concepts from text-to-image diffusion models. Most of them (Gandikota et al.



Figure 4: The comparison of the effects when unlearning multiple concepts

2023; Kumari et al. 2023; Gandikota et al. 2024; Zhang et al. 2024; Zhao et al. 2024) focus on the generation module, the U-Net. Specifically, (Gandikota et al. 2023; Kumari et al. 2023; Zhao et al. 2024) focus on updating the entire U-Net parameters, while (Gandikota et al. 2024; Zhang et al. 2024) focus on updating the cross-attention part of the U-Net parameters. (Fuchi and Takagi 2024) focus on updating the parameters of the text encoder. Updating the U-Net parameters comes with high computational complexity and is time-consuming. Additionally, updating the U-Net may influence the delicate generation and further affect the fidelity of the generated image. In contrast, updating the text encoder does not affect the fidelity and is time-efficient. Since some text-to-image and text-to-video models utilize a shared text encoder (Wang et al. 2023; Rombach et al. 2022), there is great potential to transfer the unlearning effect from image domain to the video domain.

Currently, there is a scarcity of research on unlearning concepts from text-to-video diffusion models. In our method, we leverage the transfer capability of the shared text encoder within text-to-video models.

Unlearning in Transformer-based Models

Several methods implement the unlearning in transformed-based models. (Chen and Yang 2023) introduced unlearn layers into the large language models to achieve unlearning. (Tian et al. 2024) utilize gradient information to address excessive unlearning in large language model.

Compare to (Chen and Yang 2023)’s work, Our method does not introduce auxiliary parameters into the unlearning model. The method proposed by (Tian et al. 2024) to overcome excessive unlearning is promising for the future work on our method.

Methodology

We aim to prevent the generation of specific concepts in video diffusion models. To achieve this, we propose transferring the unlearning capability of text-to-image diffusion models to text-to-video models. Specifically, we select video and image diffusion models that share a common text encoder. We then implement unlearning on the text encoder in the text-to-image domain to transfer the unlearned ability to the text-to-video domain. Importantly, we fix the U-Nets of both models, ensuring that the models’ generation capability is preserved.

Our method builds upon (Fuchi and Takagi 2024)’s work. In their work, the author only update the parameters of the text encoder, and they lists the reasons. First, research suggests that the quality of text-image alignment correlates with the quality of the text encoder (Saharia et al. 2022). DALLE-3 (Betker et al. 2023) successfully achieves high-quality generation by choosing GPT4 (Achiam et al. 2023) as its text encoder. Second, the output of text encoder is a multi-dimensional vector that preserves meaningful information about the concepts described in the text.

In addition to the benefit discussed above, the transferability of the text encoder is also important. Text encoders trained on the text-to-image and text-to-video domains may play different roles in semantic embedding. Specifically, one role is for image semantic embedding, and the other is for video semantic embedding. However, we postulate that unlearning concepts in the text encoder can have transfer capability to both its text-to-image and text-to-video domain semantic embedding.

Preliminary of Diffusion methods

Diffusion models (Ho, Jain, and Abbeel 2020; Song, Meng, and Ermon 2020) are proposed to learn the underlying distribution of the data through diffusion and denoising process. For input data $\mathbf{x} \sim p(\mathbf{x})$, in order to construct noisy sample $\mathbf{x}_t = \alpha_t \mathbf{x} + \sigma_t \epsilon$, the diffusion process adds random noise $\epsilon \sim \mathcal{N}(0, \mathbf{I})$ to it to form a Markov Chain of T steps. The noise scheduler is parameterized by the diffusion time step t , α_t and σ_t . Specifically, $SNR = \log(\alpha_t^2 / \sigma_t^2)$, the signal-to-noise ratio monotonically decreases over time. The models gradually denoise a normal distributed variable to learn the reverse process of the fixed Markov Chain of T steps to optimize a variational lower bound on $p(\mathbf{x})$. These models are weighted sequence of denoising U-Nets $\epsilon_\theta(\mathbf{x}, t)$ which predict the noises of \mathbf{x}_t . The learning objective is as follows:

$$L_{DM} = \mathbb{E}_{\mathbf{x}, \epsilon \sim \mathcal{N}(0, \mathbf{I}), t} \|\epsilon - \epsilon_\theta(\mathbf{x}_t, t)\|. \quad (1)$$

Latent diffusion models (Rombach et al. 2022) use a variational autoencoder structure. The encoder \mathcal{E} compresses the input data into low-dimensional latent representation $\mathbf{z} = \mathcal{E}(\mathbf{x})$. Unlike direct diffusion models, the diffusion and denoising processes of latent diffusion models are implemented in the latent space. This setting saves substantial training and inference time. In the final denoising stage, the output is decoded as $\mathcal{D}(\mathbf{z}_0)$ which is the reconstructed data.

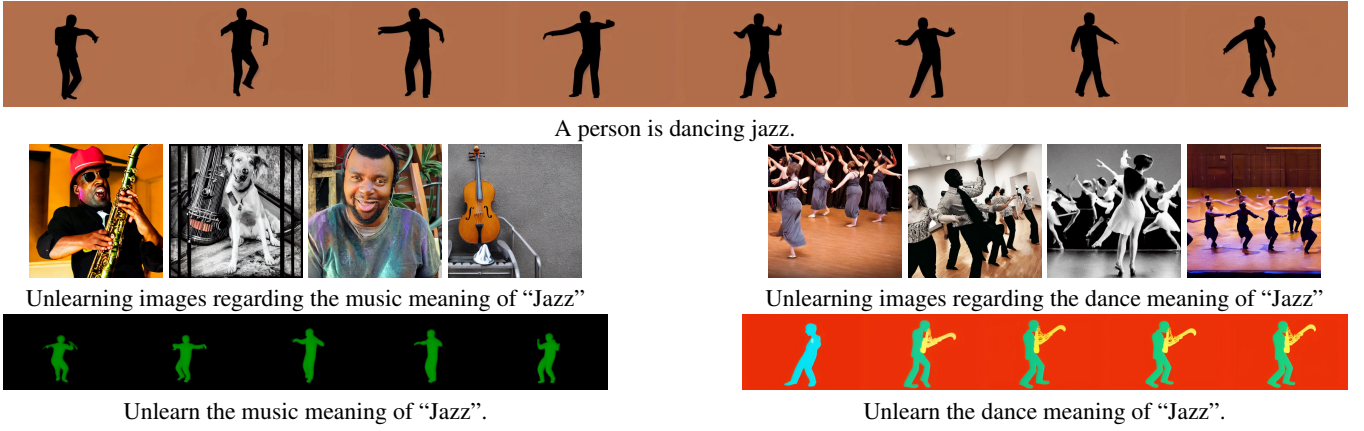


Figure 5: The comparison of unlearning the different meaning of polysemous concepts.

The objective of latent diffusion models is as follows:

$$L_{LDM} = \mathbb{E}_{\mathcal{E}(\mathbf{x}), \mathbf{e} \sim \mathcal{N}(0, \mathbf{I}), t} \|\mathbf{e} - \mathbf{e}_{\theta}(\mathbf{z}_t, t)\|. \quad (2)$$

The loss of latent diffusion models conditioned on the text input \mathbf{y} is as follows:

$$L_{LDM} = \mathbb{E}_{\mathcal{E}(\mathbf{x}), \mathbf{y}, \mathbf{e} \sim \mathcal{N}(0, \mathbf{I}), t} \|\mathbf{e} - \mathbf{e}_{\theta}(\mathbf{z}_t, t, \mathbf{c}_{\beta}(\mathbf{y}))\| \quad (3)$$

where \mathbf{c}_{β} is the text encoder output and β is the parameters of the text encoder.

Unlearning method

Because we are going to unlearn the concept in the text-to-image domain and transfer that unlearning capability to the text-to-video domain. The symbol \mathbf{x} discussed below represents images rather than videos, although videos are also feasible with significantly greater resource consumption in the following optimization.

In order to guarantee the semantic meaning of most of the concepts while changing only a small number of concepts, we apply a slight change to the text encoder's \mathbf{c}_{β} parameters

$$\mathbf{c}_{\beta} \leftarrow \mathbf{c}_{\beta} + \delta \beta. \quad (4)$$

In order to unlearn a specific concept \mathbf{y}^* , a common unlearning method is to implement gradient ascent with respect to the parameters that can be optimized. The \mathbf{x}^* represent the images described the concept \mathbf{y}^* . According to (Fuchi and Takagi 2024), we optimize the β , the parameters of the text encoder, with respect to the loss

$$L_{LDM} = \mathbb{E}_{\mathcal{E}(\mathbf{x}^*), \mathbf{y}, \mathbf{e} \sim \mathcal{N}(0, \mathbf{I}), t} \|\mathbf{e} - \mathbf{e}_{\theta}(\mathbf{z}_t, t, \mathbf{c}_{\beta}(\mathbf{y}^*))\| \quad (5)$$

while keeping the other parameters frozen.

In order to make only slight changes to variables and minimize the impact on other concepts, we implement gradient ascent for only 5 epochs with respect to \mathbf{x}^* . \mathbf{x}^* is practically represented by 4 images to facilitate few-shot learning. These images are generated by text-to-image diffusion models to reduce the workload of collection.

After unlearning, we use the same text encoder, but we apply it for text-to-video generation. The procedure of our method is shown in Figure 2.

Experiment

We conduct qualitative experiments and ablation studies. Since there currently exists no other method for unlearning concepts in text-to-video diffusion models, we do not present a comparison of our method with other methods.

Experiment Setup

We apply our model on Lavie(Wang et al. 2023) by transferring the text encoder, which is unlearned on the stable diffusion 1.5(Rombach et al. 2022). We use four generated images in the few shot setting to unlearn a specific concept. We optimize the text encoder using Adam(Kingma and Ba 2014). The hyperparameters used in the experiment are listed in Table 1.

Table 1: Hyperparameters

Hyperparameter	Value
Training epochs	5
Batch size	2
Learning rate	10^{-5}
Weight decay	10^{-8}
Adam(β_1, β_2)	(0.9, 0.98)

Qualitative Results

We analyzed the qualitative results of unlearning a single object, its effect on other concepts, unlearning multiple concepts and unlearning concepts with multiple meanings.

Unlearning Single Object Single-object unlearning experiments are illustrated in Figure 1. We conducted the experiments on the concepts including ‘‘Snoopy’’, ‘‘Van Gogh’’, ‘‘Eiffel Tower’’ and ‘‘Elon Musk’’. Our method successfully replaced the Snoopy with a shadow of a wolf walking in the park at sunrise. Furthermore, our method removed the Van Gogh style orange background of a panda taking a selfie and replaced it with a realistic green background. Additionally, for videos with the prompt ‘‘a butterfly flying on the Eiffel Tower’’, our method removed the Eiffel Tower throughout the videos, leaving only a single butterfly flying in the

frames. Finally, for the prompt “Elon Musk in a space suit standing besides a rocket, high quality”, our method removed the Elon Musk’s facial features, protecting his privacy. These experiments demonstrated that our method can unlearn copyrighted cartoon characters, artist styles, objects and people’s facial characteristics.

Effect on Other Concepts We conducted an experiment to unlearn the concept “Eiffel Tower”. We compared the influence on the generation results of three other prompts. The results were shown in the Figure 3. There was no significant effect on the generations of “a polar bear playing drum kit in NYC Times Square, 4k, high resolution” and “jungle river at sunset, ultra quality”. However, there seems to be an effect on the generation of the concept “Tokyo Tower” in the prompt “a cat climbing on Tokyo Tower”. In the original generation, there was a Tokyo Tower in the scene. In the generation after unlearning “Eiffel Tower”, the Tokyo Tower seems to be interpreted as a general building. This experiment show that the unlearning procedure only influences the video generation of the concepts similar to target concepts and guarantees the generation quality of other concepts.

Unlearning Multiple Concepts As people may be interested in unlearning multiple copyright concepts within a single text encoder, we conducted experiments focused on unlearning “Snoopy” and “Mickey Mouse.” In cases without unlearning, the successful generations showed “a Snoopy walking in the park at sunrise, oil painting style” and “Mickey Mouse is greeting the children.” When unlearning “Snoopy,” the first prompt’s generation replaced Snoopy with the shadow of a wolf, while the second prompt successfully generated an image without Mickey Mouse’s influence. When unlearning both “Snoopy” and “Mickey Mouse,” the first prompt’s generation again omitted Snoopy, and the second prompt’s generation excluded Mickey Mouse. These experiments demonstrate the feasibility of unlearning multiple copyrighted concepts simultaneously.

Unlearning Concepts with Multiple Meanings Since some concepts have multiple meanings and people want to unlearn a specific meaning of those concepts, we conducted experiments on unlearning “Jazz”. “Jazz” has two meanings. One correlates with the music and is represented by the saxophone. The other meaning is a dancing style. We aimed to unlearn the music meaning of the “Jazz” by using several images related to saxophones and music performance. Similarly, we unlearn the dance meaning of “Jazz” by using several images related to dance performance. To test the unlearning effect in the text-to-video domain, we used the prompt “A person is dancing jazz.”. The generation results of the unlearning process are shown in the Figure 5. When unlearning the music meaning of “Jazz”, the video still showed the person dancing jazz. Conversely, when unlearning the dance meaning of “Jazz”, the videos showed the person playing saxophones, which is the characteristic of the music meaning of “Jazz”. These experiments demonstrate the possibility of unlearning specific meanings of the polysemous concepts which are copyright-related.

Because there is no proper metric for evaluating unlearn-

ing performance, we did not implement quantitative experiments.

Ablation Studies

Our method involves the k -shot unlearning process in the text-to-image text encoder and the training epoch number of the unlearning process. In the following experiment, we demonstrate the influence of different settings for transfer learning on the text-to-video generation.

k -shot Unlearning When implementing the unlearning, we need to collect several concept-related images to help unlearn the concept of the text encoder and then to transfer the unlearning effect to the text-to-video diffusion models. Normally, we use the four generated images as the concept-related images and it is the four-shot learning. In this experiment, we compare zero-shot, two-shot and four-shot cases.

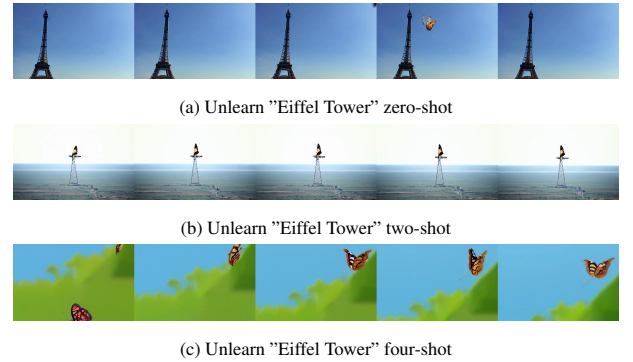


Figure 6: The comparison of the effects after unlearning “Eiffel Tower” in different shots. The prompt is “a butterfly flying on the Eiffel Tower”.

As shown in the Figure 6, the concepts and images of Eiffel Tower were gradually removed as the number of shots increased. During unlearning “Eiffel Tower” with zero-shot images, the Eiffel Tower concepts and images remained preserved in the video. When unlearning “Eiffel Tower” with two-shot images, the video depicted the Eiffel Tower as a high-voltage power tower. It seems that the high-voltage power tower is similar to the Eiffel Tower. Finally, with four-shot images, the Eiffel Tower was completely removed from the video.

Number of Epochs To implement unlearning, we need to train for several epochs to unlearn the concepts learned by the text encoder. Then, we transfer this unlearning effect to the text-to-video diffusion models. Typically, five epochs are used for training. In this experiment, we compare the training results obtained with one, two, three, four, and five epochs.

As we can see in Figure 7, with an increasing number of epochs, Elon Musk’s facial characteristics disappear. When the epoch number is 1, the generated man’s facial characteristics are similar to Elon Musk. After epoch number 2, Elon Musk’s facial characteristics disappear, and we cannot tell the identity of the generated man.



Figure 7: The comparison of the effects after unlearning “Elon Musk” in different epoch trainings. The prompt is “Elon Musk in a space suit standing besides a rocket, high quality”.

Limitations

The experiments demonstrate the effectiveness of our transfer method. However, based on previous experiments, we observed that when the model unlearns a specific concept, the generation of similar concepts can be influenced. For instance, when unlearning “Eiffel Tower,” the generation of “Tokyo Tower” was affected, and the semantic meaning of “Tokyo Tower” seemed to be mapped to “building.” Furthermore, we found that when a specific concept is unlearned, the new meaning associated with the concept name may not be significantly similar to its original semantic meaning. This contradicts the conclusion presented in (Fuchi and Takagi 2024).

Conclusion

Our method achieves unlearning concepts by transferring the unlearning capability of the text encoder from text-to-image diffusion models to text-to-video diffusion models. The process of unlearning a specific concept involves optimizing the parameters of the text encoder through gradient ascent on the objective function based on several concept-related images. The optimized text encoder is then reused in the text-to-video diffusion models. Since the unlearning procedure focuses on the text-to-image domain and only optimizes the text encoder, it is fast, taking only about 100s to unlearn a concept on an RTX 3070. In our experiment, we found that several target concepts disappeared in the videos while the effect on other, different concepts was minimal. The models can also jointly unlearn multiple concepts. Additionally, our ablation study demonstrate the need for a sufficient number of images for few-shot learning and for the

models to be optimized by a sufficient number of epochs.

Our method can unlearn static concepts like objects, artistic styles, cartoon characters, and human appearances. However, it cannot handle dynamic concepts such as specific dance routines, acrobatics, and continuous photographic works. These are difficult to express in the text-to-image domain and are also protected by copyright law. We plan to investigate how to unlearn dynamic concepts in the future. This may involve directly optimizing the text encoder for the text-to-video domain. Ultimately, we aim to combine optimizations in both text-to-image and text-to-video domains to achieve unlearning of dynamic concepts.

References

- Achiam, J.; Adler, S.; Agarwal, S.; Ahmad, L.; Akkaya, I.; Aleman, F. L.; Almeida, D.; Altschmidt, J.; Altman, S.; Anadkat, S.; et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Betker, J.; Goh, G.; Jing, L.; Brooks, T.; Wang, J.; Li, L.; Ouyang, L.; Zhuang, J.; Lee, J.; Guo, Y.; et al. 2023. Improving image generation with better captions. *Computer Science*. <https://cdn.openai.com/papers/dall-e-3.pdf>, 2(3): 8.
- Bourtole, L.; Chandrasekaran, V.; Choquette-Choo, C. A.; Jia, H.; Travers, A.; Zhang, B.; Lie, D.; and Papernot, N. 2021. Machine unlearning. In *2021 IEEE Symposium on Security and Privacy (SP)*, 141–159. IEEE.
- Chen, J.; and Yang, D. 2023. Unlearn what you want to forget: Efficient unlearning for llms. *arXiv preprint arXiv:2310.20150*.
- Fuchi, M.; and Takagi, T. 2024. Erasing Concepts from Text-to-Image Diffusion Models with Few-shot Unlearning. *arXiv preprint arXiv:2405.07288*.
- Gandikota, R.; Materzynska, J.; Fiotto-Kaufman, J.; and Bau, D. 2023. Erasing concepts from diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2426–2436.
- Gandikota, R.; Orgad, H.; Belinkov, Y.; Materzyńska, J.; and Bau, D. 2024. Unified concept editing in diffusion models. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 5111–5120.
- Goodfellow, I.; Pouget-Abadie, J.; Mirza, M.; Xu, B.; Warde-Farley, D.; Ozair, S.; Courville, A.; and Bengio, Y. 2014. Generative adversarial nets. *Advances in neural information processing systems*, 27.
- Goodfellow, I.; Pouget-Abadie, J.; Mirza, M.; Xu, B.; Warde-Farley, D.; Ozair, S.; Courville, A.; and Bengio, Y. 2020. Generative adversarial networks. *Communications of the ACM*, 63(11): 139–144.
- Gu, X.; Du, C.; Pang, T.; Li, C.; Lin, M.; and Wang, Y. 2023. On memorization in diffusion models. *arXiv preprint arXiv:2310.02664*.
- Ho, J.; Chan, W.; Saharia, C.; Whang, J.; Gao, R.; Gritsenko, A.; Kingma, D. P.; Poole, B.; Norouzi, M.; Fleet, D. J.; et al. 2022. Imagen video: High definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303*.

- Ho, J.; Jain, A.; and Abbeel, P. 2020. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33: 6840–6851.
- Ho, J.; and Salimans, T. 2021. Classifier-Free Diffusion Guidance. In *NeurIPS 2021 Workshop on Deep Generative Models and Downstream Applications*.
- Kingma, D. P.; and Ba, J. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Kumari, N.; Zhang, B.; Wang, S.-Y.; Shechtman, E.; Zhang, R.; and Zhu, J.-Y. 2023. Ablating concepts in text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 22691–22702.
- Ma, R.; Zhou, Q.; Xiao, B.; Jin, Y.; Zhou, D.; Li, X.; Singh, A.; Qu, Y.; Keutzer, K.; Xie, X.; et al. 2024. A Dataset and Benchmark for Copyright Protection from Text-to-Image Diffusion Models. *arXiv preprint arXiv:2403.12052*.
- Ni, J.; Abrego, G. H.; Constant, N.; Ma, J.; Hall, K. B.; Cer, D.; and Yang, Y. 2021. Sentence-t5: Scalable sentence encoders from pre-trained text-to-text models. *arXiv preprint arXiv:2108.08877*.
- Radford, A.; Kim, J. W.; Hallacy, C.; Ramesh, A.; Goh, G.; Agarwal, S.; Sastry, G.; Askell, A.; Mishkin, P.; Clark, J.; et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, 8748–8763. PMLR.
- Raffel, C.; Shazeer, N.; Roberts, A.; Lee, K.; Narang, S.; Matena, M.; Zhou, Y.; Li, W.; and Liu, P. J. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140): 1–67.
- Rombach, R.; Blattmann, A.; Lorenz, D.; Esser, P.; and Ommer, B. 2022. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 10684–10695.
- Ronneberger, O.; Fischer, P.; and Brox, T. 2015. 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*, 234–241. Springer.
- Saharia, C.; Chan, W.; Saxena, S.; Li, L.; Whang, J.; Denton, E. L.; Ghasemipour, K.; Gontijo Lopes, R.; Karagol Ayan, B.; Salimans, T.; et al. 2022. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in neural information processing systems*, 35: 36479–36494.
- Setty, R. 2023. Ai art generators hit with copyright suit over artists’ images. *Bloomberg Law*. Accessed on February, 1: 2023.
- Song, J.; Meng, C.; and Ermon, S. 2020. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*.
- Tian, B.; Liang, X.; Cheng, S.; Liu, Q.; Wang, M.; Sui, D.; Chen, X.; Chen, H.; and Zhang, N. 2024. To Forget or Not? Towards Practical Knowledge Unlearning for Large Language Models. *arXiv preprint arXiv:2407.01920*.
- Wang, Y.; Chen, X.; Ma, X.; Zhou, S.; Huang, Z.; Wang, Y.; Yang, C.; He, Y.; Yu, J.; Yang, P.; et al. 2023. Lavie: High-quality video generation with cascaded latent diffusion models. *arXiv preprint arXiv:2309.15103*.
- Yin, S.; Wu, C.; Yang, H.; Wang, J.; Wang, X.; Ni, M.; Yang, Z.; Li, L.; Liu, S.; Yang, F.; et al. 2023. Nuwa-xl: Diffusion over diffusion for extremely long video generation. *arXiv preprint arXiv:2303.12346*.
- Zhang, G.; Wang, K.; Xu, X.; Wang, Z.; and Shi, H. 2024. Forget-me-not: Learning to forget in text-to-image diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 1755–1764.
- Zhao, M.; Zhang, L.; Zheng, T.; Kong, Y.; and Yin, B. 2024. Separable Multi-Concept Erasure from Diffusion Models. *arXiv preprint arXiv:2402.05947*.