

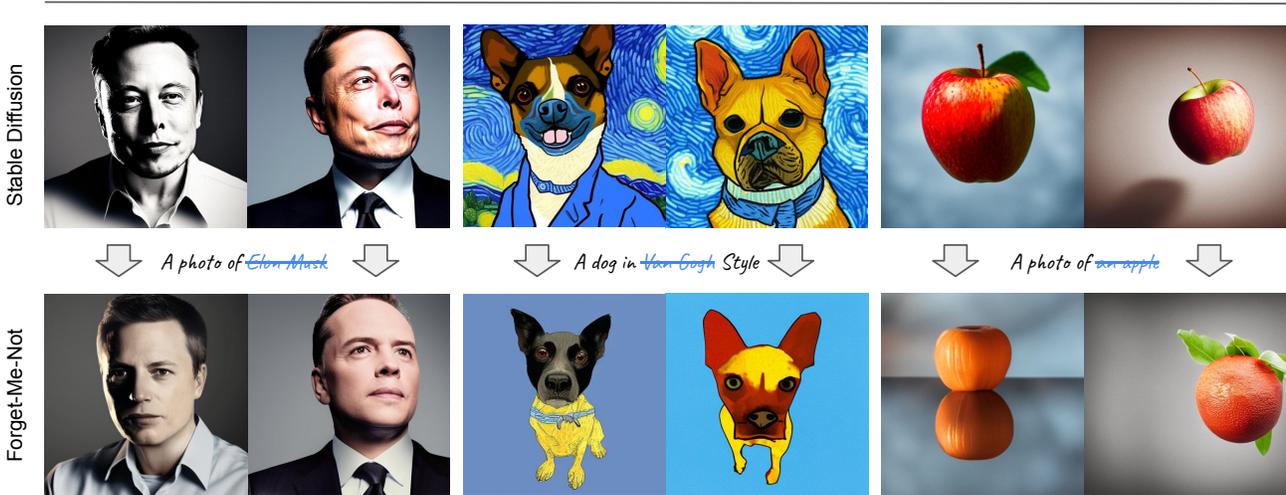
Forget-Me-Not: Learning to Forget in Text-to-Image Diffusion Models

Eric Zhang^{1*}, Kai Wang^{1*}, Xingqian Xu^{1,3}, Zhangyang Wang^{2,3}, Humphrey Shi^{1,3}

¹SHI Labs @ U of Oregon & UIUC, ²UT Austin, ³Picsart AI Research (PAIR)

<https://github.com/SHI-Labs/Forget-Me-Not>

Concept Forgetting



Concept Correction & Disentangle



Figure 1: Given a text-to-image model (*i.e.* Stable Diffusion), our approach can swiftly re-steer the cross attention towards a specific concept and subsequently forgetting or correcting it. (1) Concept Forgetting: target concepts (denoted in blue text and crossed-out) are successfully removed without compromising the quality of the output. (2) Concept Correction & Disentangle: our method can be used to correct a dominant or undesired concept of a prompt. Prior overshadowed concepts reveal in outputs after the dominant concepts are forgotten. In addition, our method learns to forget fast with only 30 seconds for certain concepts (e.g. Elon Musk), and can be easily adapted to lightweight model patches for Stable Diffusion, allowing for multi-concept manipulation and convenient distribution to users.

Abstract

The unlearning problem of deep learning models, once primarily an academic concern, has become a prevalent issue in the industry. The significant advances in text-to-image generation techniques have prompted global discus-

sions on privacy, copyright, and safety, as numerous unauthorized personal IDs, content, artistic creations, and potentially harmful materials have been learned by these models and later utilized to generate and distribute uncontrolled content. To address this challenge, we propose **Forget-Me-Not**, an efficient and low-cost solution designed to safely remove specified IDs, objects, or styles from a well-configured

*Equal contribution

text-to-image model in as little as 30 seconds, without impairing its ability to generate other content. Alongside our method, we introduce the **Memorization Score (M-Score)** and **ConceptBench** to measure the models' capacity to generate general concepts, grouped into three primary categories: ID, object, and style. Using M-Score and ConceptBench, we demonstrate that Forget-Me-Not can effectively eliminate targeted concepts while maintaining the model's performance on other concepts. Furthermore, Forget-Me-Not offers two practical extensions: a) removal of potentially harmful or NSFW content, and b) enhancement of model accuracy, inclusion and diversity through **concept correction and disentanglement**. It can also be adapted as a lightweight model patch for Stable Diffusion, allowing for concept manipulation and convenient distribution. To encourage future research in this critical area and promote the development of safe and inclusive generative models, we will open-source our code and ConceptBench at <https://github.com/SHI-Labs/Forget-Me-Not>.

1. Introduction

Recently, text-to-image models [10, 23, 53, 52, 57, 66, 54, 65] have shown impressive performance in synthesizing high-quality images according to text prompts. Among these methods, diffusion models such as DALL-E 2 [52] and Stable Diffusion (SD) [54] have met commercial-level productization requirements, initiating numerous applications for downstream users; such text-to-images are also recently shown to be able to generate and editing videos in a zero-shot fashion [33] without further training. Industrial solutions such as [49, 21, 35, 48, 44, 60] have been widely adopted in various art and visual design systems, garnering significant public attention. Despite the popularity of this field, concerns about security, fairness, regulation, copyright, safety, etc., continue to grow rapidly in proportion to model usages. Risks such as generating unauthorized, biased, and unsafe content have become an immediate issue to be resolved. While this is not the first time the community has investigated these cases, prior efforts such as [22, 32, 4, 36] have proposed methods in which most of them are high-cost solutions focused on GAN. Yet we still need an effective and efficient solution that can be widely applied to large-scale diffusion models, which motivated us to dive deep into this topic.

The risks and issues associated with such large-scale text-to-image models originate from the billion-sized datasets used in training, including public datasets such as Laion [58], COYO [5], CC12M [11], and private data from Google [57, 66], OpenAI [53, 52], etc. The public datasets are usually web-scraped images and captions that lack human-level quality assurance on bias and safety, while private data sources are impossible to determine at scale. As

a result, it is nearly unfeasible to fully address harmful content, privacy and copyright concerns through data filtering or source attribution. A compromised solution could be domain adaptation [25, 67, 64]. In practice, people can adapt a large-scale model to a clean small/mid-size dataset and later use the model for image synthesis. However, collecting and filtering such datasets may still be quite laborious. Worse than that, such domain adaptation has severely influenced model capacity, making out-of-domain image synthesis challenging and sometimes nearly impossible.

Will there be another path? Designing efficient methods and algorithms that guide existing large-scale text-to-image models to *forget certain concepts* could be a better solution. We start this paper by first introducing this new mission, namely *concept forgetting*, in which a designated set of concepts can be safely disentangled from the visual content. To achieve this goal, we proposed **Forget-Me-Not**, a simple, low-cost, but effective solution for concept forgetting. We also proposed the **memorization score (M-score)** along with **ConceptBench**, in which the former gauge the generative power of models on certain concepts, and the latter introduce sets of benchmark to assess concept forgetting and memorization. Moreover, we extend concept forgetting to *concept correction & disentangle* that may further assist models in being accurate and diverse.

In conclusion, the main contribution of this paper can be summarized as the following:

- We propose **Forget-Me-Not**, a plug-and-play, efficient and effective concept forgetting and correction method for large-scale text-to-image models. It provides an efficient way to forget specific concepts with as few as 35 optimization steps, which typically takes about 30 seconds. Additionally, Forget-Me-Not can be easily adapted as lightweight patches for Stable Diffusion, allowing for multi-concept manipulation and convenient distribution to text-to-image model users to address privacy, copyright, and safety concerns.
- We also propose **memorization score (M-score)** and **ConceptBench**, which enable quantitative measurements of models' capacity for synthesizing the target set of concepts. To the best of our knowledge, we are the first to introduce a numerical solution to gauge the model's behavior of memorization and forgetting.
- Through extensive studies and tests, we demonstrate that our Forget-Me-Not is simple, low-cost, and effective. Downstream applications, such as harmful and NSFW content removal and biased concept correction & disentangle, further expand our scope beyond concept forgetting towards cross-modal model refinements that may better fit real-world use cases.

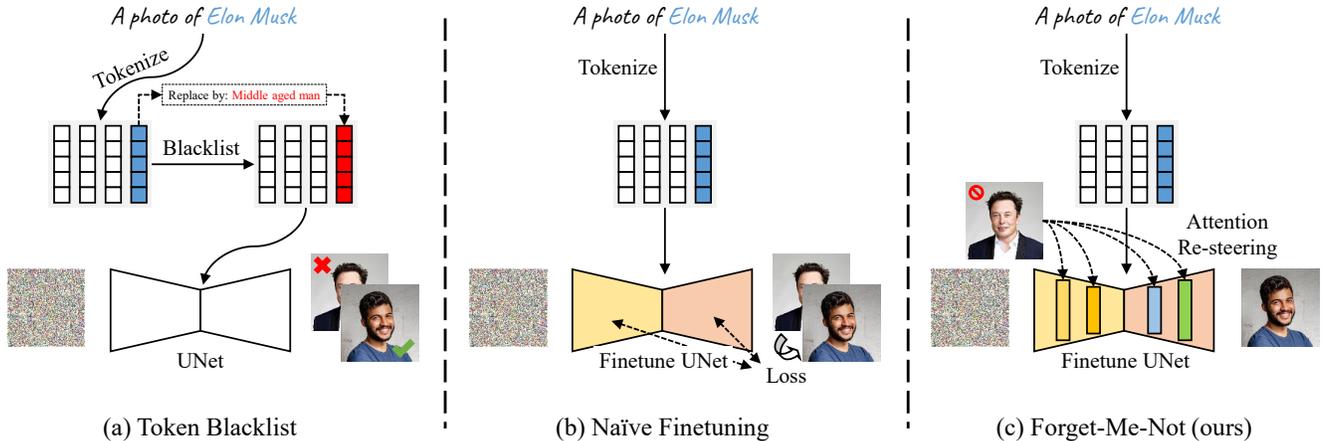


Figure 2: This figure shows two baseline forgetting methods and our proposed Forget-Me-Not. The target concept to forget is Elon Musk. One baseline is (a) Token Blacklist that simply replaces the target token with a different one. The other baseline is (b) Naïve Finetuning in which instead of replacing tokens, it finetunes model weights so that the new weights generate outputs containing unrelated concepts. Our method (c) Forget-Me-Not utilizes Attention Re-steering in which we finetune only UNet to minimize each of the intermediate attention maps associated with the target concepts to forget.

2. Related Works

2.1. Text-to-Image Synthesis

Image generation has been a challenging but very attractive research area, whose goal is to synthesize natural-looking images. In the past decade, we have witnessed the rapid advance of it from unconditional generative models to conditional generative models with powerful architectures of auto-regressive model [53, 66], GAN [8, 31, 30, 63, 59] and diffusion process [26, 47, 41, 19, 3, 61]. Early works focus on unconditional, single-category data distribution modeling, such as hand-written digits, certain species of animals, and human faces [16, 12, 31, 40]. Though, unconditional models quickly achieves photo realistic results among single-category data, it’s shown that mode collapsing issue usually happens when extending data distributions to multiple-category or real image diversity [8, 43, 1].

To tackle the model collapsing problem, the conditional generative model has been introduced. Since then, different types of data have been used as the conditioning for generative models, e.g. class labels, image instances, and even networks [8, 45] etc. At the same time, CLIP [50, 28], a large-scale pretrained image-text contrastive model, provides a text-image prior of extremely high diversity, which is discovered to be applicable as the conditioning for generative model [46, 15, 38]. Nowadays, DALL-E 2 [52] and Stable Diffusion [54] are capable of generating high quality images solely conditioning on free-form texts, inheriting the diversity of billions of real images scraping from the Internet. Subsequently, a line of work seeks to efficiently adapt the massive generative model to generate novel rendition of an unseen concept represented by a small reference

set, leveraging the great diversity. Dreambooth [56] proposed to adapt the model by finetuning all of its weights, while it requires enormous storage to save newly adapted weights. Textual Inversion [24] and LoRA [27] ameliorate the issue by adapting the model by adding a small set of extra weights.

2.2. Model Unlearning

However, this great diversity comes at a price. It incurs potential risk of privacy leakage and copyright infringement. [7, 60] have successfully retrieved samples from Stable Diffusion that are highly faithful to real training examples. Therefore, being able to forget/unlearn certain concept in a model without hurting the generative ability for the rest is of both research and practical interests. Similar topics have been seen in fields other than conditional generative modeling. In model-agnostic meta-learning, [2] noted selectively forgetting the influence of prior knowledge in a network improves the performance in adapted tasks. [9, 39, 42, 13] explores the unlearning of a set of requested data points in a pretrained model.

Our work differs from existing forgetting and unlearning works in a few aspects. First, we study forgetting in the context of text-to-image generative models. Second, we are deleting not only requested data points represented by a small reference set, but the concept behind those data points, which possesses significant impact in text-to-image generation due to the fact that it’s almost impossible to enumerate all prompts and synonyms relating to a concept.

3. Method

3.1. Preliminaries

Diffusion models [26, 47, 18] are denoising models that iteratively restore data x_0 from its Gaussian noise corruption x_T with a total step number T . Such a restoration process is usually known as the reverse diffusion process $p_\theta(x_{t-1}|x_t)$ and the opposite of the reverse process is the forward diffusion process that blends the signal with noise $q(x_t|x_{t-1})$:

$$q(x_t|x_{t-1}) = \mathcal{N}(x_t; \sqrt{1 - \beta_t}x_{t-1}; \beta_t\mathbf{I})$$

$$p_\theta(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t); \Sigma_\theta(x_t, t))$$

Both forward and reverse processes are presumably Markovian chains, so we can express the likelihood of both processes as:

$$q(x_{1:T}|x_0) = \prod_{t=1}^T q(x_t|x_{t-1})$$

$$p_\theta(x_{0:T}) = p(x_T) \prod_{t=1}^T p_\theta(x_{t-1}|x_t)$$

The loss function for the diffusion process is then to minimize the variational bound \mathcal{L}_{vlb} of the negative log-likelihood $p_\theta(x_0)$ (i.e. maximize the likelihood of x_0 as the final denoised result from a model with parameters θ):

$$\begin{aligned} \mathcal{L}_{VLB} &= \mathbb{E}[-\log p_\theta(x_0)] \\ &\leq \mathbb{E}_q \left[-\log \frac{p_\theta(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] \end{aligned}$$

Cross-Attentions [62] are widely adopted deep learning modules used in discriminative models [20, 6, 29], conditional generation models [52, 57, 54] as well as language models [17, 51, 14]. The purpose of cross-attention is to transfer information from conditional inputs to hidden features through dot product and softmax. For example, in stable diffusion [54], the hidden feature serves as the query Q and context serves as key K and value V . Assume Q and K has dimension d for inner product, the output h is then computed as the following:

$$h = \text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right)V$$

It is important to note that such QKV assignments are not fixed. Other assignments, such as conditional-driven queries and feature-driven keys and values, may also have their usage.

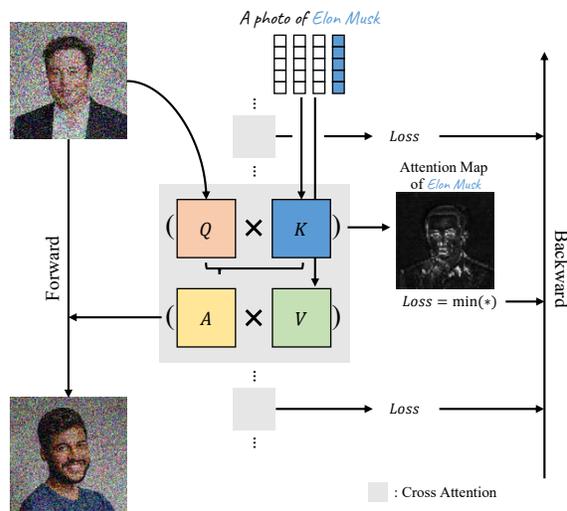


Figure 3: This figure shows the Attention Re-steering we proposed in our Forget-Me-Not method, in which we set the objective function to minimize the attention maps of target concepts (i.e. Elon Musk in this case) and correspondingly finetune the network.

3.2. Concept Forgetting

A concept is an abstract term representing an intuited object of thought, which also serves as the foundation for people’s perceptions. Specifically for computer vision, we may recognize concepts as tangible things, including identities, objects that physically existed, style of images, object relations, and even poses and behavior. Concept forgetting, literarily speaking, is the action of reverting a model from understanding certain concepts. On contrast to machine unlearning, which aims to delete the fields around designated data points, we define concept forgetting in diffusion models as the disentanglement of concept prompts and visual contents. This definition allows models to retain their generative abilities to the greatest extent possible.

Besides, we set the following four goals for concept forgetting research:

- **Performance:** the proposed approach should at best remove target concepts from the model.
- **Integrity:** the proposed approach should at best keep other concepts of the model.
- **Generality:** the proposed approach can be applied to a wide range of concepts that covers all aspects of human perceptions.
- **Flexibility:** the proposed approach can be applied to various models of different tasks and domains.

3.3. Forget-Me-Not

To fulfill the end goals we mentioned in Section 3.2, we introduce **Forget-Me-Not**, a heuristically and important ap-

Methods	Performance	Integrity	Generality	Flexibility
Token Blacklisting	No forgetting	Inevitably affects other concepts sharing overlapping prompts	Within the vocabulary of the tokenizer	Tokenizer required
Naive Finetuning	Successfully removes concept	Removes unrelated concept by fault	Applies to any concepts with sufficient data.	Applies to any models
Forget-Me-Not	Successfully removes concept	Maintains most of the model’s integrity.	Applies to any concepts with few data samples	Only applies to models with cross attention

Table 1: This table compares pros (green) and cons (red) on the four major aspects of concept forgetting between baselines and the proposed Forget-Me-Not. If an approach can handle an aspect to some extent, the corresponding explanation is marked in yellow.

belong to the more general “dog” category but have distinct visual features. In the case of brands, they represent abstract concepts of intangible objects that can manifest as logos throughout our daily lives.

Object is a broader concept encompassing multiple variations. For example, “dog” refers to various breeds of dogs that share common features. By combining identity instances mentioned earlier, this category provides a hierarchical structure to examine the influence of concept forgetting on the model. We include food items like “apple”, “banana”, and “broccoli”, man-made objects such as “airplane”, “keyboard”, “motorcycle”, “umbrella”, and “boat”, and general animals like “dog” and “horse.”

Style is an abstract concept that determines the overall appearance of generated images. ConceptBench incorporates styles such as “Van Gogh”, “Picasso”, “doodle”, “pixel art”, “neon”, and “sketch.”

4.2. Baseline

In view of the multi-component nature of Stable Diffusion models, there are several naive methods that can be used to superficially remove a concept from them, such as blacklisting keywords in prompts, removing specific tokens from the tokenizer dictionary, or tuning the model with unrelated images to divert the target concept, as illustrated in Figure 2(a)(b). However, these methods can result in a significant deterioration and shifting of the model’s generation capability. Removing tokens from the dictionary can alter the tokenization of prompts where those tokens were previously used and affect the generation of other prompts with overlapping tokens. For instance, removing tokens of “Hillary Clinton” could lead to dysfunctionality in generating “Bill Clinton”. Naive finetuning to forget with unrelated images explicitly overwrites the visual representation of a concept with extra data and runs the risk of compromising



Figure 4: Finetuning to forget concept “Johnny Depp” with unrelated images of “a photo of man”. This method distorts other concepts with visual details of selected unrelated images.

existing concept space, as shown in Figure 4. Moreover, it is impossible to exhaust test all relation-based concepts for blacklisting or finetuning.

4.3. Qualitative Comparison

We present the results of concept forgetting from our benchmark, illustrated in Figure 5, the Multi-concepts model of Elon Musk and Taylor Swift demonstrates our method’s ability to perform multi-concept forgetting. As shown in the first row, both target concepts have been forgotten. We evaluated the impacts of forgetting specific concepts on other related concepts, examining four related concepts to Elon Musk and Taylor Swift - man, woman, Bill Gates, and Emma Watson. As shown, Forget-Me-Not achieved good content preservation and visual quality. However, we observed minor pose and style changes in man and Bill Gates. Based on these findings, we posit that our approach may have a greater impact on closely related concepts than on others. Additionally, the last row shows that a

Multi-Concept Model: Elon Musk & Taylor Swift

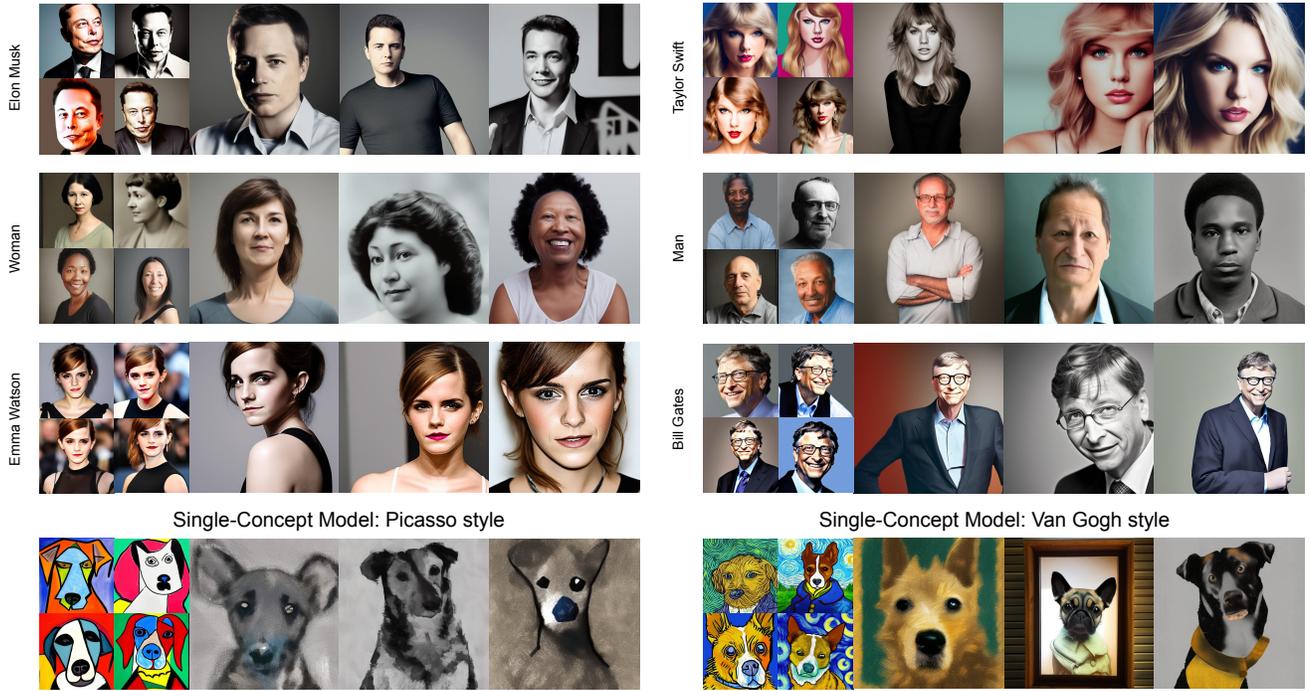


Figure 5: Results of concept forgetting using our method. The first 2x2 grid shows the original samples in Stable Diffusion. The subsequent 3 images are sampled after concept forgetting, using the same prompt. The top 3 rows are from a multi-concept model targeting both Elon Musk and Taylor Swift, demonstrating the multi-concept forgetting capability. Control concepts such as Bill Gates and Emma Watson manifest that our approach has minimal impact on concepts other than target ones. The last row shows two single-concept model of styles. Output images were generated with prompts: “a photo of X” (top 3 rows), “a dog in X style” (bottom row).



Figure 6: Concept Correction: Once the dominant concept has been diminished by our method, the lesser concepts of an semantic-rich prompt become more prominent in generated results. Output images were generated with prompts (top to bottom): “a movie poster of Mulan”, “James Bond”, “apple shape”.

new painting style is emerging after forgetting Picasso and Van Gogh styles.



Figure 7: In concept correction, our method has the advantage of comprehensive forgetting over negative prompt.

4.4. Quantitative Analysis

Memorization Measurement Textual inversion [24] can be employed to identify the token embeddings that best correspond to images. We utilize this technique to measure the concept embedding changes of anchor images toward a reference before and after forgetting. These changes can be seen as generative model’s memorization level of a concept, which we term Memorization Score.

In the case of “Elon Musk”, prompt “Elon Musk” is used as reference. Its concept embedding (\mathbf{emb}_r) is obtained by processing it through text encoder. Subsequently, we invert the same anchor images of Elon Musk using *original model* and *forgetting model* respectively. Concept embeddings of anchor images are obtained by processing inverted tokens through text encoder, resulting in two sets: original textual inversion (\mathbf{emb}_o) and forgetting textual inversion (\mathbf{emb}_f). Only the pooler tokens of concept embeddings are used for measurement. The change in concept embedding is quantified as the difference between $\cos(\mathbf{emb}_r, \mathbf{emb}_o)$ and $\cos(\mathbf{emb}_r, \mathbf{emb}_f)$. A decrease indicates successful forgetting. Since the textual inversion process bring the randomness of embeddings, we compute the average Memorization Score over five running times. We present memorization scores from each sub-category in Table 2. Additional results can be found in the Supplementary material.

4.5. Concept Correction

It has been observed that in text-to-image models, the semantics of a prompt are often dominated by the one with the most number of image-text examples in the training set, resulting in the suppression of inferior semantics during in-

Concept	Initial Mem Score	Forgetting Mem Score
Elon Musk	0.943	0.848
Mickey Mouse	0.948	0.836
Zebra	0.972	0.899
Google	0.940	0.811
Apple	0.696	0.493
Horse	0.877	0.808
Van Gogh	0.916	0.684

Table 2: Memorization Scores of instances from each sub-categories.

ference. Figure 6 exemplifies this scenario, where generation is dominated by a concept that is strongly correlated with a prompt due to unbalanced training examples. In the case of the James Bond series, the generation results are overwhelmingly dominated by Daniel Craig, as shown in the middle row. However, our method manages to diminish the most prominent semantic in the prompt, i.e., Daniel Craig, and make other James Bond actors visible. Similarly, in the case of Mulan series and the homonym of “apple”, where the apple fruit and Apple brand are competing with each other, our method successfully corrects target concept in generated images .

Negative prompt is a technique used in text-to-image synthesis to eliminate unwanted concepts associated with a prompt. However, their use can result in negative impacts on other aspects of the image, such as changes to its structure and style. Furthermore, negative prompts fail to correct undesirable concepts under certain circumstances. For example, in Figure 7, “a photo of a mango” consistently generates dog images. This is because the name “mango” is commonly used as a pet name for dogs, and people upload photos of their dogs to the internet, which are collected as training data. In this case, using a negative prompt for dogs would be ineffective, as mango is also a popular cat name, creating the problem of endlessly expanding negative prompts. However, our method successfully brings back the mango fruit by forgetting the connection between “a photo of mango” and dog/cat images.

4.6. NSFW Removal

In this section, we examine the effectiveness of our method in a real case of removing harmful contents. NSFW is an internet shorthand for “not safe for work”, used for indicating contents that are not wished to be seen in the public. Such content may include material even offensive for adult audiences. However, they inevitably present in large datasets such as LAION [58], even though NSFW detector has been used [34]. Stable Diffusion, trained on LAION, is known for generating NSFW content when prompted with certain triggers.

To evaluate our method, we use a well-known NSFW-

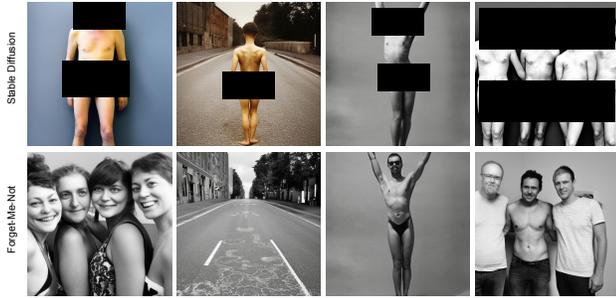


Figure 8: Results of removing NSFW contents triggered by “naked”. Faces and sensitive parts are blacked out.

triggering prompt, “a photo of naked” in Stable Diffusion v2.1 model. Using EulerAncestralDiscreteScheduler, inference-step 50, and scale-guidance 8, the model consistently generates images containing nude individuals. We use eight generated NSFW images as input for training Forget-Me-Not.

The results, shown in Figure 8, indicate that the concept of “naked” has been successfully forgotten. Notably, the second row shows that all sensitive visual cues have been changed in different ways. The first example changes abruptly from a naked man to a group of smiling women. In the second example, NSFW individual has been removed from the scene. The last two examples render clothed people, making them safe. Overall, our method achieved efficient forgetting of NSFW content without the need for additional data or the assistance of third-party NSFW detectors.

4.7. Ablation Studies

Concept Inversion Ablation We conducted experiments on concepts from ConceptBench, with and without using concept inversion (CI). Concept inversion is used to handle concepts that are difficult to describe using prompts. Generally, it can help extract the target concept from the prompt, resulting in more precise embeddings. However, precise embeddings may not be always ideal, see Section 4.7 Concept Ablation. Our results show that CI can achieve higher fidelity for concepts that can be well-described in a prompt, as illustrated in Figure 9, where the model trained with CI preserved more of the original poses and details.

Trainable Weights Ablation We conducted experiments to compare finetuning the entire UNet model versus only finetuning the cross-attention (CA) layers. Cross-attention is a critical component in text-to-image generation, as it injects textual information into the image formation process. Given the same hyper-parameter settings except for steps, our results show that both methods can successfully achieve concept forgetting. However, finetuning the entire UNet model tended to break the model’s generation capability in fewer steps. In some cases, the model collapsed before the forgetting process was complete, as show



Figure 9: Improving fidelity to original model with concept inversion. Concept prompt tend to have diverse semantics, resulting in distortion in concept forgetting. CI extracts precise semantics into dedicated tokens, allowing for more pose and feature consistency.

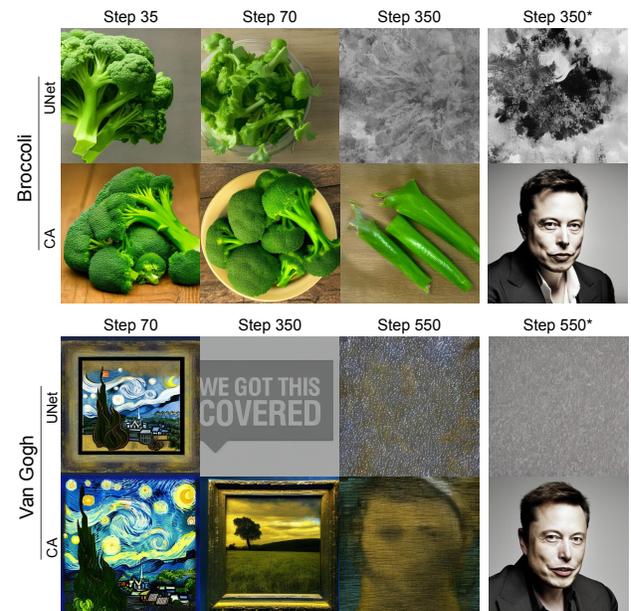


Figure 10: Trainable weights ablation using UNet and Cross Attention (CA). Compared to CA, UNet is more sensitive to optimization steps. The last column with Step X^* shows the control concept Elon Musk at Step X .

in the “Broccoli” case of Figure 10.

Token Embedding of Concept Ablation Our method relies on token embeddings of a concept, which are critical for computing the attention re-steering loss. As shown in Section 4.7 on Concept Inversion Ablation, changing the token embeddings of a concept produces varying results. In Figure 11, we demonstrate a situation where concept prompt prevails concept inversion. By using the same set-



Figure 11: Comparison of different token embeddings in concept forgetting. Given concept airplane, we compare forgetting with either tokens of “airplane” or inverted tokens using concept inversion, where forgetting with CI fails.

tings except for token embeddings, prompt of “airplane” succeeds while inverted tokens fails. We hypothesize that minimizing cross attention over these specific inverted tokens of “airplane” tends to break generative capability of the model quickly.

5. Conclusion

In this study, we investigate concept forgetting in text-to-image generative models and introduce Forget-Me-Not. This lightweight approach enables ad-hoc concept forgetting using only a few either real or generated concept images; it can also be easily distributed using model patches. Forget-Me-Not is further naturally extended to enable concept correction and disentanglement. Our experiments demonstrate that Forget-Me-Not is successful in diminishing and correcting target concepts in Stable Diffusion. Additionally, we introduce ConceptBench and Memorization Score as evaluation metrics. Overall, our work provides a foundation for further research on concept forgetting and manipulation in text-to-image generation, and can be further extended to other conditional multimodal generative models to improve the accuracy, inclusion and diversity of such models.

6. Social Impact & Limitations

Social Impact Our research has a positive social impact by offering an effective and cost-efficient method to remove and correct harmful and biased concepts in text-to-image generative models. These models are rapidly becoming the backbone of popular AI art and graphic design tools, used by a growing number of people. Our method can generate lightweight model patches that can be conveniently distributed to text-to-image model users like how conventional software patch works. Thus, our research takes a small

step towards promoting fairness and privacy protection in AI tools, ultimately benefiting society as a whole.

Limitations While our approach performs well on concrete concepts in ConceptBench, it faces challenges in identifying and forgetting abstract concepts. Additionally, successful forgetting may require manual interventions, such as concept-specific hyperparameter tuning.

References

- [1] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks. In *International Conference on Machine Learning (ICML)*, pages 214–223. PMLR, 2017.
- [2] Sungyong Baik, Seokil Hong, and Kyoung Mu Lee. Learning to forget for meta-learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2379–2387, 2020.
- [3] Yogesh Balaji, Seungjun Nah, Xun Huang, Arash Vahdat, Jiaming Song, Karsten Kreis, Miika Aittala, Timo Aila, Samuli Laine, Bryan Catanzaro, et al. ediffi: Text-to-image diffusion models with an ensemble of expert denoisers. *arXiv preprint arXiv:2211.01324*, 2022.
- [4] Hugo Berg, Siobhan Mackenzie Hall, Yash Bhargat, Wonsuk Yang, Hannah Rose Kirk, Aleksandar Shtedritski, and Max Bain. A prompt array keeps the bias away: Debiasing vision-language models with adversarial learning. *arXiv preprint arXiv:2203.11933*, 2022.
- [5] Minwoo Byeon, Beomhee Park, Haecheon Kim, Sungjun Lee, Woonhyuk Baek, and Saehoon Kim. Coyo-700m: Image-text pair dataset. <https://github.com/kakaobrain/coyo-dataset>, 2022.
- [6] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *European Conference on Computer Vision (ECCV)*, pages 213–229. Springer, 2020.
- [7] Nicholas Carlini, Jamie Hayes, Milad Nasr, Matthew Jagielski, Vikash Sehwal, Florian Tramèr, Borja Balle, Daphne Ippolito, and Eric Wallace. Extracting training data from diffusion models. *arXiv preprint arXiv:2301.13188*, 2023.
- [8] Arantxa Casanova, Marlene Careil, Jakob Verbeek, Michal Drozdal, and Adriana Romero Soriano. Instance-conditioned gan. *Advances in Neural Information Processing Systems (NeurIPS)*, 34:27517–27529, 2021.
- [9] Sungmin Cha, Sungjun Cho, Dasol Hwang, Honglak Lee, Taesup Moon, and Moontae Lee. Learning to unlearn: Instance-wise unlearning for pre-trained classifiers. *arXiv preprint arXiv:2301.11578*, 2023.
- [10] Huiwen Chang, Han Zhang, Jarred Barber, AJ Maschinot, Jose Lezama, Lu Jiang, Ming-Hsuan Yang, Kevin Murphy, William T Freeman, Michael Rubinstein, et al. Muse: Text-to-image generation via masked generative transformers. *arXiv preprint arXiv:2301.00704*, 2023.
- [11] Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. Conceptual 12m: Pushing web-scale image-text pre-training to recognize long-tail visual concepts. In *Proceed-*

- ings of the *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3558–3568, 2021.
- [12] Yunjey Choi, Youngjung Uh, Jaejun Yoo, and Jung-Woo Ha. Stargan v2: Diverse image synthesis for multiple domains. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8188–8197, 2020.
- [13] Vikram S Chundawat, Ayush K Tarun, Murari Mandal, and Mohan Kankanhalli. Zero-shot machine unlearning. *arXiv preprint arXiv:2201.05629*, 2022.
- [14] Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*, 2022.
- [15] Katherine Crowson, Stella Biderman, Daniel Kornis, Dashiell Stander, Eric Hallahan, Louis Castricato, and Edward Raff. Vqgan-clip: Open domain image generation and editing with natural language guidance. In *European Conference on Computer Vision (ECCV)*, pages 88–105. Springer, 2022.
- [16] Li Deng. The mnist database of handwritten digit images for machine learning research [best of the web]. *IEEE Signal Processing Magazine*, 29(6):141–142, 2012.
- [17] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [18] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems (NeurIPS)*, 34:8780–8794, 2021.
- [19] Tim Dockhorn, Arash Vahdat, and Karsten Kreis. Genie: Higher-order denoising diffusion solvers. *arXiv preprint arXiv:2210.05475*, 2022.
- [20] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [21] Adobe Firefly. <https://www.adobe.com/sensei/generative-ai/firefly.html>.
- [22] Felix Friedrich, Patrick Schramowski, Manuel Brack, Lukas Struppek, Dominik Hintersdorf, Sasha Luccioni, and Kristian Kersting. Fair diffusion: Instructing text-to-image generation models on fairness. *arXiv preprint arXiv:2302.10893*, 2023.
- [23] Oran Gafni, Adam Polyak, Oron Ashual, Shelly Sheynin, Devi Parikh, and Yaniv Taigman. Make-a-scene: Scene-based text-to-image generation with human priors. In *European Conference on Computer Vision (ECCV)*, pages 89–106. Springer, 2022.
- [24] Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or. An image is worth one word: Personalizing text-to-image generation using textual inversion. *arXiv preprint arXiv:2208.01618*, 2022.
- [25] Giorgio Giannone, Didrik Nielsen, and Ole Winther. Few-shot diffusion models. *arXiv preprint arXiv:2205.15463*, 2022.
- [26] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems (NeurIPS)*, 33:6840–6851, 2020.
- [27] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- [28] Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori, Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hananeh Hajishirzi, Ali Farhadi, and Ludwig Schmidt. OpenCLIP, 7 2021.
- [29] Jitesh Jain, Jiachen Li, MangTik Chiu, Ali Hassani, Nikita Orlov, and Humphrey Shi. Onerformer: One transformer to rule universal image segmentation. *arXiv preprint arXiv:2211.06220*, 2022.
- [30] Tero Karras, Miika Aittala, Samuli Laine, Erik Härkönen, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Alias-free generative adversarial networks. *Advances in Neural Information Processing Systems (NeurIPS)*, 34:852–863, 2021.
- [31] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4401–4410, 2019.
- [32] Patrik Joslin Kenfack, Daniil Dmitrievich Arapov, Rasheed Hussain, SM Ahsan Kazmi, and Adil Khan. On the fairness of generative adversarial networks (gans). In *2021 International Conference "Nonlinearity, Information and Robotics"(NIR)*, pages 1–7. IEEE, 2021.
- [33] Levon Khachatryan, Andranik Movsisyan, Vahram Tadevosyan, Roberto Henschel, Zhangyang Wang, Shant Navasardyan, and Humphrey Shi. Text2video-zero: Text-to-image diffusion models are zero-shot video generators. *arXiv preprint arXiv:2303.13439*, 2023.
- [34] LAION-AI. Clip-based-nsfw-detector. <https://github.com/-AI/CLIP-based-NSFW-Detector>.
- [35] Lexica. <https://lexica.art/>.
- [36] Zhiheng Li, Anthony Hoogs, and Chenliang Xu. Discover and mitigate unknown biases with debiasing alternate networks. In *European Conference on Computer Vision (ECCV)*, pages 270–288. Springer, 2022.
- [37] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European Conference on Computer Vision (ECCV)*, pages 740–755. Springer, 2014.
- [38] Xihui Liu, Dong Huk Park, Samaneh Azadi, Gong Zhang, Arman Chopikyan, Yuxiao Hu, Humphrey Shi, Anna Rohrbach, and Trevor Darrell. More control for free! image synthesis with semantic diffusion guidance. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 289–299, 2023.
- [39] Yang Liu, Zhuo Ma, Ximeng Liu, Jian Liu, Zhongyuan Jiang, Jianfeng Ma, Philip Yu, and Kui Ren. Learn to for-

- get: Machine unlearning via neuron masking. *arXiv preprint arXiv:2003.10933*, 2020.
- [40] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Large-scale celebfaces attributes (celeba) dataset. *Retrieved August, 15(2018):11*, 2018.
- [41] 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. *arXiv preprint arXiv:2206.00927*, 2022.
- [42] Ananth Mahadevan and Michael Mathioudakis. Certifiable machine unlearning for linear models. *arXiv preprint arXiv:2106.15093*, 2021.
- [43] Luke Metz, Ben Poole, David Pfau, and Jascha Sohl-Dickstein. Unrolled generative adversarial networks. *arXiv preprint arXiv:1611.02163*, 2016.
- [44] Midjourney. <https://www.midjourney.com/>.
- [45] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014.
- [46] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021.
- [47] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *International Conference on Machine Learning (ICML)*, pages 8162–8171. PMLR, 2021.
- [48] NovelAI. <https://novelai.net/>.
- [49] Picsart. <https://picsart.com/>.
- [50] 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 (ICML)*, pages 8748–8763. PMLR, 2021.
- [51] 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. *The Journal of Machine Learning Research*, 21(1):5485–5551, 2020.
- [52] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.
- [53] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine Learning (ICML)*, pages 8821–8831. PMLR, 2021.
- [54] 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 (CVPR)*, pages 10684–10695, June 2022.
- [55] 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.
- [56] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. 2022.
- [57] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S Sara Mahdavi, Rapha Gontijo Lopes, et al. Photorealistic text-to-image diffusion models with deep language understanding. *arXiv preprint arXiv:2205.11487*, 2022.
- [58] 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. *arXiv preprint arXiv:2210.08402*, 2022.
- [59] Tamar Rott Shaham, Tali Dekel, and Tomer Michaeli. Singan: Learning a generative model from a single natural image. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 4570–4580, 2019.
- [60] Gowthami Somepalli, Vasu Singla, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Diffusion art or digital forgery? investigating data replication in diffusion models. *arXiv preprint arXiv:2212.03860*, 2022.
- [61] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020.
- [62] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in Neural Information Processing Systems (NeurIPS)*, 30, 2017.
- [63] Steven Walton, Ali Hassani, Xingqian Xu, Zhangyang Wang, and Humphrey Shi. Stylenat: Giving each head a new perspective. *arXiv preprint arXiv:2211.05770*, 2022.
- [64] Jiayu Xiao, Liang Li, Chaofei Wang, Zheng-Jun Zha, and Qingming Huang. Few shot generative model adaption via relaxed spatial structural alignment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11204–11213, 2022.
- [65] Xingqian Xu, Zhangyang Wang, Eric Zhang, Kai Wang, and Humphrey Shi. Versatile diffusion: Text, images and variations all in one diffusion model. *arXiv preprint arXiv:2211.08332*, 2022.
- [66] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Guncan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, et al. Scaling autoregressive models for content-rich text-to-image generation. *arXiv preprint arXiv:2206.10789*, 2022.
- [67] Jingyuan Zhu, Huimin Ma, Jiansheng Chen, and Jian Yuan. Few-shot image generation with diffusion models. *arXiv preprint arXiv:2211.03264*, 2022.