

# Nepotistically Trained Generative-AI Models Collapse

Matyáš Boháček  
Stanford University  
maty@stanford.edu

Hany Farid  
University of California, Berkeley  
hfarid@berkeley.edu

## Abstract

Trained on massive amounts of human-generated content, AI-generated image synthesis is capable of reproducing semantically coherent images that match the visual appearance of its training data. We show that when retrained on even small amounts of their own creation, these generative-AI models produce highly distorted images. We also show that this distortion extends beyond the text prompts used in retraining, and that once poisoned, the models struggle to fully heal even after retraining on only real images.

## 1 Introduction

From text to audio and image, today’s generative-AI systems are trained on large quantities of human-generated content. Most of this content is obtained by scraping a variety of online sources. As generative AI becomes more common, it is reasonable to expect that future data scraping will invariably catch generative AI’s own creations. We ask what happens when these generative systems are trained on varying combinations of human-generated and AI-generated content.

Although it is early in the evolution of generative AI, there is already some evidence that retraining a generative AI model on its own creation – what we call model poisoning – leads to a range of artifacts in the output of the newly trained model. It has been shown, for example, that when retrained on their own output, large language models (LLMs) contain irreversible defects that cause the model to produce gibberish – so-called model collapse [22].

Similarly, on the image generation side, it has been shown [1] that when retrained on its own creations, StyleGAN2 [12] generates images (of faces or digits) with visual and structural defects. Interestingly, the authors found that there was a deleterious effect on output as the ratio of AI-generated content used to retrain the model ranged from 0.3% to 100%.

It has also been shown that in addition to GAN-based image generation, diffusion-based text-to-image models are also vulnerable. The authors in [15, 14], for example, showed that in a simplified setting, retraining on one’s own creation can lead to



**Figure 1:** Representative examples of images generated by the baseline version of Stable Diffusion (prompt: “older hispanic man”).

image degradation and a loss of diversity. Examined the impact of poisoning of ID-DPM [16], the authors in [8] report somewhat conflicting results depending on the task at hand (recognition, captioning, or generation). With respect to generation, the authors report a relatively small impact on image quality but do note a lack of diversity in the retrained models.

Building on these earlier studies, we show that the popular open-source model Stable Diffusion (SD)<sup>1</sup> is highly vulnerable to data poisoning. In particular, we show that when iteratively retrained on faces of its own creation the model – after an initial small improvement – quickly collapses, yielding highly distorted and less diverse faces. Somewhat surprisingly, even when the retraining data contains only 3% of self-generated images, this model collapse persists. We also investigate the extent of this model poisoning beyond the prompts used for retraining, and examine the ability of the poisoned model to heal with further retraining on only real images.

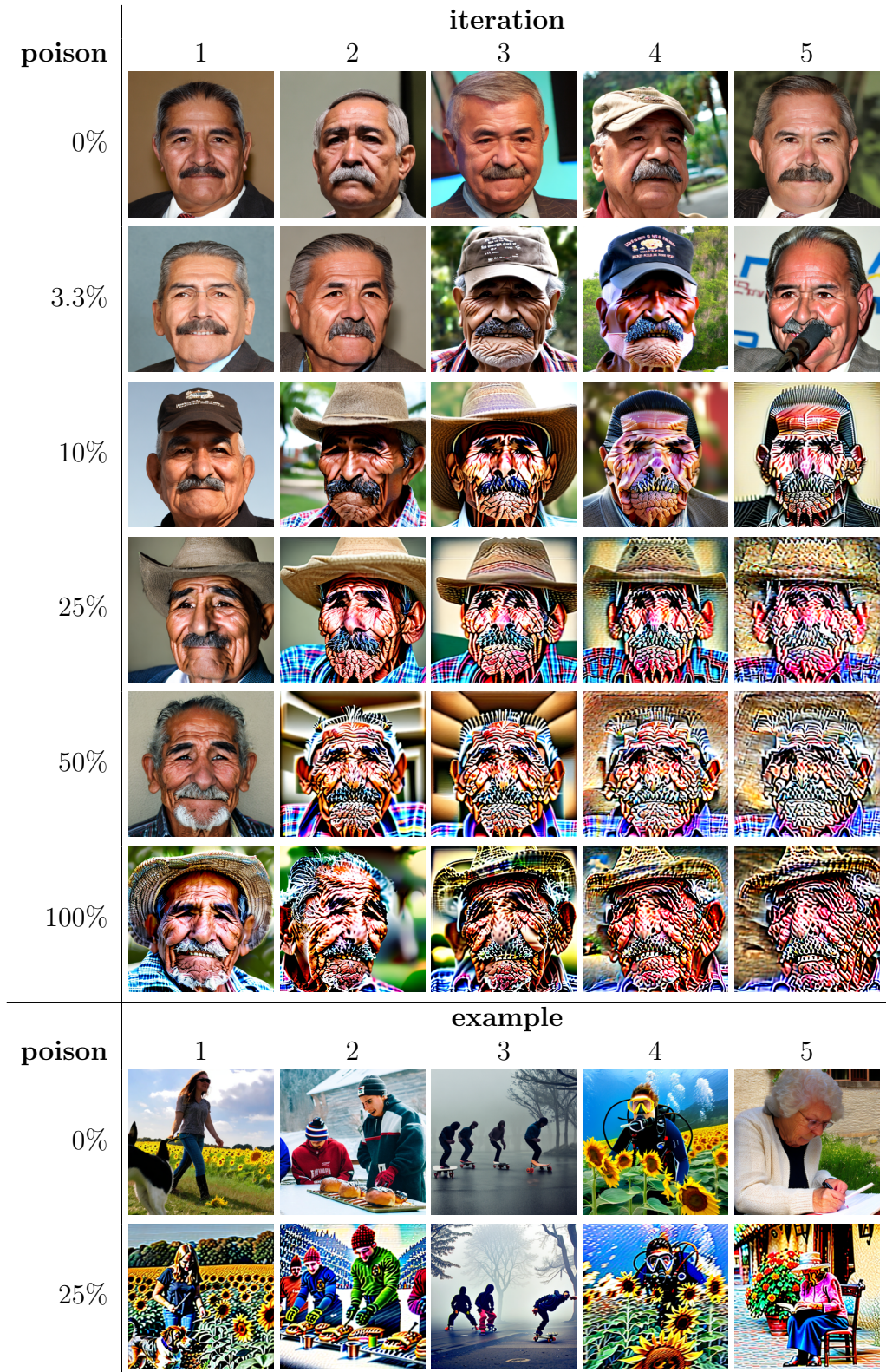
## 2 Results

Shown in Figure 1 are five representative images generated from the baseline Stable Diffusion (SD) model for a single demographic group (“older hispanic man”). Generally speaking, images from the baseline model are consistent with the text prompt and of high visual quality.

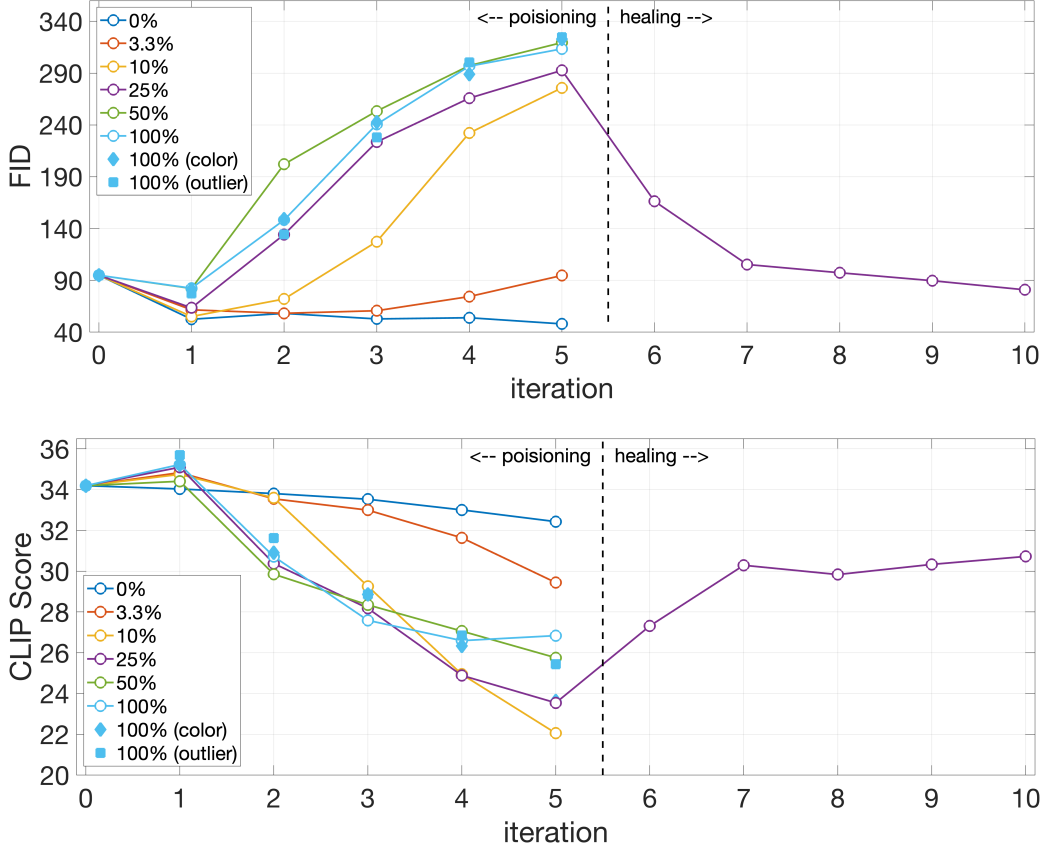
Shown in the first row of Figure 2 are representative images generated from iterative retraining (Section 3.2) of the baseline SD model on images of real faces taken from the FFHQ facial dataset (Section 3.1). The generated images are semantically consistent with the text prompts, exhibit the prototypical alignment property of the faces in the FFHQ dataset, and show no signs of distortion. Shown in Figure 3 are the FID and CLIP score (Section 3.3) for the full set of generated images (labeled 0%), from which we see that iterative retraining on real images causes no degradation in the resulting model.

Also shown in top portion of Figure 2 are representative images generated from iterative retraining the baseline SD model with a different mixture of self-generated and real images. Regardless of the mixture ratio, the iterative retraining eventually leads to collapse by the fifth iteration at which point the generated images are highly distorted. This collapse is quantified in Figure 3 where we see that after a small improvement in quality on the first iteration, both the FID and CLIP score

<sup>1</sup><https://github.com/Stability-AI/StableDiffusion>



**Figure 2:** Representative examples generated after iterative retraining for different compositions of the retraining dataset ranging from (top to bottom) 0% SD-generated and 100% real to 100% SD-generated faces and 0% real faces. Shown in the lower panel are representative images generated with text prompts distinct from those used in the model retraining.



**Figure 3:** Shown is the Fréchet inception distance (FID) and contrastive language-image pre-training score (CLIP) as a function of the number of retraining iterations and the composition of the retraining dataset ranging from 100% SD-generated faces and 0% real faces to 0% SD-generated and 100% real (“poisoning”). The diamond plot symbol corresponds to the 100%/0% condition in which the retraining dataset is color matched to the real faces. The square plot symbol corresponds to the the same condition in which the retraining dataset was curated on each iteration to remove low quality faces. The trend is the same for both metrics: the presence of generated faces leads to a degradation in quality across iterations (a higher FID and a lower CLIP correspond to lower image quality). See also Figure 2. Also shown is the FID and CLIP score for the 25% model as it is retrained for another five iterations on only real images (“healing”).

(Section 3.3) reveal a significant degradation in image quality (a high FID and a low CLIP correspond to lower quality images).

The open plot symbols (diamond and square) in Figure 2 correspond to the two control conditions (Section 3.4) in which the retraining dataset is color matched to real images (diamond) and any low-quality generated images are replaced with high-quality generated images prior to retraining (square). Even these curated retraining datasets lead to model collapse at the same rate as the other datasets.

In addition to the degradation in image quality, and consistent with previous reports [8], we also note that model poisoning leads to a lack of diversity in terms of

the appearance of the generated faces. This can be seen in Figure 2 where, particularly when the poisoning is larger than 10%, the generated faces are highly similar across the latter iterations.

Shown in the lower two rows of Figure 2 are representative examples of images generated with text prompts<sup>2</sup> distinct from the demographic prompts used in the model retraining (see Section 3.1). As expected, the images generated by the model retrained on entirely real images (0%) produces semantically coherent images with no obvious visual artifacts. However, the images generated by the model retrained on 25% SD-generated faces often – but not always – exhibits the same textured artifacts as seen in the faces in the upper portion of this figure. This means that the model poisoning is not limited to a specific category of images used in the retraining, but seems to impact the model more broadly.

Lastly, we wondered if, once poisoned, the model could be “healed” by retraining on only real images. The model poisoned for five iterations with 25% SD-generated images was retrained for another five iterations on only real images. Shown in Figure 4 are representative examples of faces generated from five different demographic groups across these five additional iterations. Although in some cases by the tenth iteration, the generated faces have fewer artifacts, in other cases, the artifacts persist. Shown in the right portion of Figure 3 are the FID and CLIP scores for these healing iterations in which we see that the FID recovers to the original base model and the CLIP score almost recovers to base model levels.

Although the mean FID and CLIP score recovers, we clearly see remnants of the poisoning in some of the faces in Figure 4. This larger variation is evident in the standard deviation of the CLIP score which is 2.8 for the base model (with a mean of 35.1) but is 4.2 for the healed model after five iterations (with a mean of 27.3). It appears that the model can partially – but not entirely – heal.

## 3 Methods

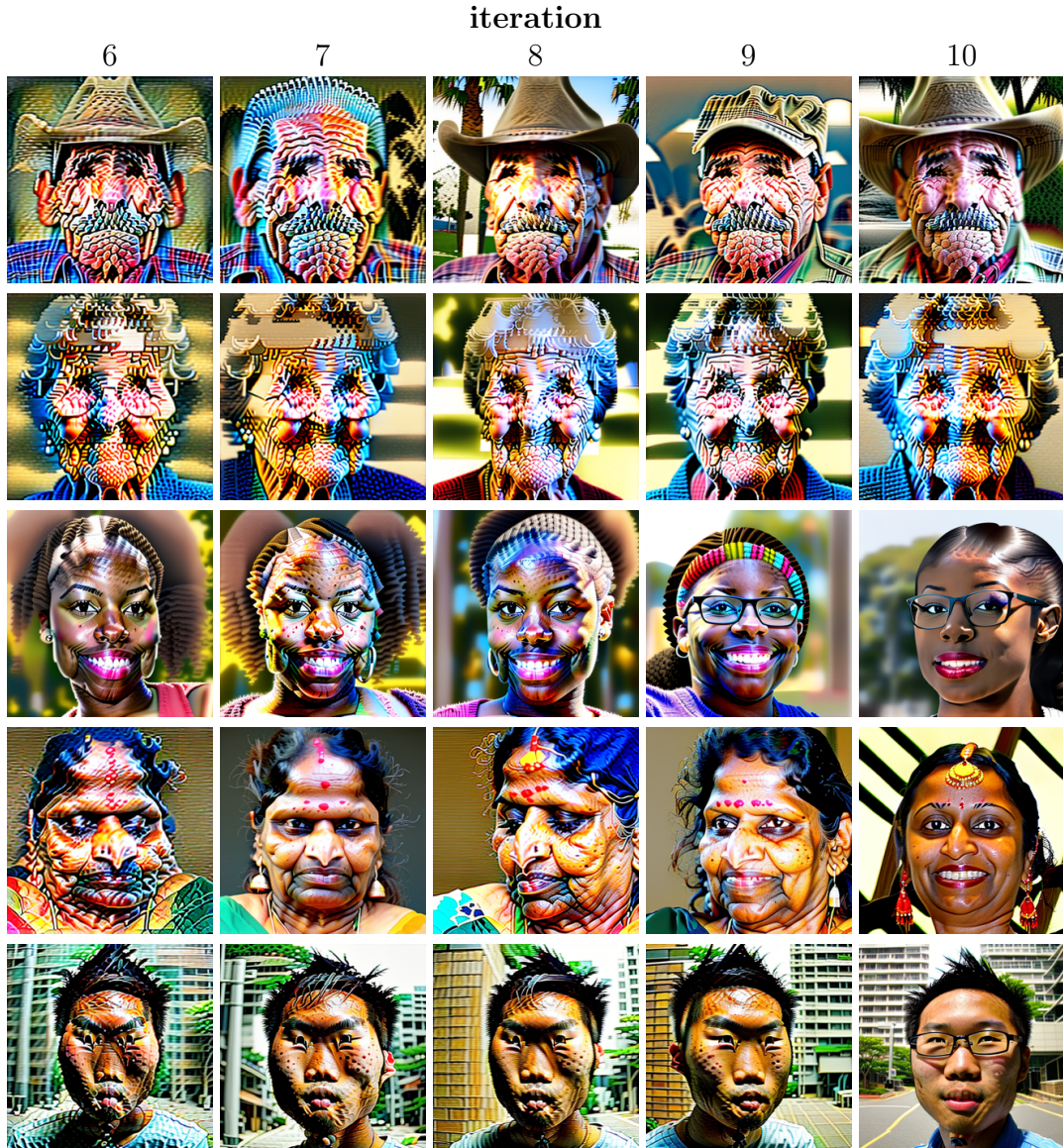
### 3.1 Images

Starting with the FFHQ image dataset containing 70,000 high-quality faces of size  $1024 \times 1024$  pixels [11], we automatically classified [20, 18] each face based on gender (man/woman), race (asian, black, hispanic, indian, and white), and age (young, middle-age, old). A total of 900 images, 30 from each of 30 (2 [gender]  $\times$  5 [race]  $\times$  3 [age]) demographic categories, were randomly selected.

We then used these real images as input to the image-to-image synthesis pipeline of Stable Diffusion (SD, v.2.1) [19] to generate 900 images consistent with the demographic prompt “a photo of a <age> <race> <gender>.” Shown in Figure 5

---

<sup>2</sup>The text prompts used to generate the images in the lower panel of Figure 2 are: “A dog walker training a dog amidst a field of sunflowers”; “A football team grilling hamburgers at a snowy ski resort”; “A group of skateboarders practicing martial arts in a mysterious foggy landscape”; “A marine biologist scuba diving amidst a field of sunflowers”; and “An elderly woman writing in a journal in a charming village square”.



**Figure 4:** Representative examples generated after iterative retraining of the 25% poisoned model with only real images. In some cases, the poisoning persists (top), while in others, the model appears to at least partially heal itself.

are representative examples of real (top) and generated faces (bottom). These 900 generated faces are, as described next, used to seed the iterative model retraining.

We used this image-to-image synthesis instead of the unconstrained text-to-image generation to balance the comparison of model retraining on healthy data and poisoned data.

### 3.2 Retraining

Leaving the CLIP text encoder and variational autoencoder (VAE) modules intact, we retrained the denoising U-Net module of the base Stable Diffusion model (SD



**Figure 5:** Representative examples of real images (top) used as a seed for image-to-image generation (bottom).

v.2.1) using the recommended parameters<sup>3</sup> (a constant learning rate of  $2 \times 10^{-6}$ , 50 epochs,  $512 \times 512$  image resolution, no mixed precision, and random horizontal flip augmentation). The model was initially retrained on the 900 SD-generated images and demographic captions described in Section 3.1. Another 900 images with the same demographic prompts were generated from this retrained model. These images were then used to retrain the model. This process was repeated for a total of five iterations.

This entire process was repeated with different compositions of faces ranging from the above 100% SD-generated and 0% real faces to a 50%/50%, 25%/75%, 10%/90%, 3.3%/96.7%, or 0%/100% (the odd-ball 3.3% composition corresponds to one of the 30 images per demographic group being SD-generated with the other 29 real).

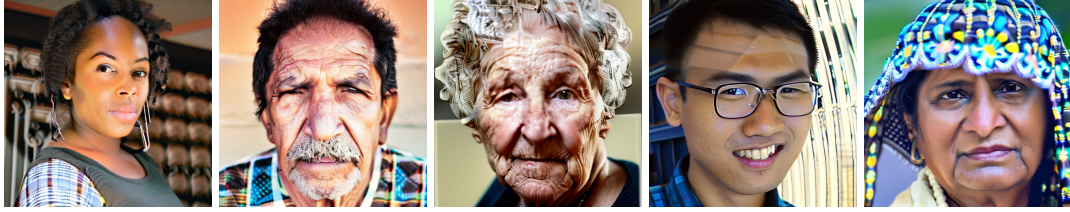
### 3.3 Evaluation

A standard Fréchet inception distance (FID) [10] and Contrastive Language-Image Pre-training score (CLIP) [17, 9] are used to assess the quality of synthesized images.

The FID compares two sets of  $N$  images with a 1 : 1 correspondence, where a smaller distance corresponds to a higher image quality. The reference images consist of 820 real faces from the FFHQ dataset (see Section 3.1). This is less than the full set of 900 because the FFHQ dataset is not sufficiently diverse across all demographic categories. We compare 820 generated images, consisting of the maximum number of images per 30 demographic groups, to this set of 820 demographically similar real images.

The CLIP score, evaluated across all 900 generated images, calculates the cosine similarity between the visual CLIP embedding of a synthesized image and the textual CLIP embedding of its caption; a larger score corresponds to higher quality. This evaluation captures the semantic consistency of the image to the input text prompt.

<sup>3</sup>[https://github.com/huggingface/diffusers/blob/ece55227ffc60f7bb09442b390a02e45ac0438f8/examples/text\\_to\\_image/train\\_text\\_to\\_image.py](https://github.com/huggingface/diffusers/blob/ece55227ffc60f7bb09442b390a02e45ac0438f8/examples/text_to_image/train_text_to_image.py)



**Figure 6:** Representative examples of low-quality generated images that are replaced in the retraining control experiment.

### 3.4 Controls

We carried out two control experiments to determine the impact to our iterative retraining of the SD model (Section 3.2).

Having noticed that, as compared to real images, SD-generated images tend to be of higher contrast and saturation, we wondered if these color differences would impact the iterative retraining. In this first control, the color histogram of each generated image is matched to a real image. Each generated image is histogram matched – in the three-channel luminance/chrominance (YCbCr) space – to a real image with the most similar color distribution (measured as the image with the minimal Kullback-Leibler divergence averaged across all three YCbCr channels). This histogram matching is performed on each retraining iteration.

We also noticed that among the 900 generated images (Section 3.1) there are occasional images with obvious artifacts, including misshapen facial features, Figure 6. In this second control, therefore, we removed from the retraining dataset any generated image with a single-image FID score [21] greater than the mean single-image FID between the first batch of 900 generated images and their corresponding 900 real images. This culling is performed on each retraining iteration.

## 4 Discussion

We find that at least one popular diffusion-based, text-to-image generative-AI system is surprisingly vulnerable to data poisoning with its own creations. This data poisoning can occur unintentionally by, for example, indiscriminately scraping and ingesting images from online sources. Or, it can occur from an adversarial attack where websites are intentionally populated with poisoned data, as described in [2]. Even more aggressive adversarial attacks can be launched by manipulating both the image data and text prompt on as little as 0.01% to 0.0001% of the dataset [3].

In the face of these vulnerabilities, there are some reasonable measures that could be taken to mitigate these risks. First, there is a large body of literature for classifying images as real or AI-generated (e.g., [24, 4, 7, 13, 5]). An ensemble of these types of detectors could be deployed to exclude AI-generated images from being ingested into a model’s retraining. A complementary approach can automatically and robustly watermark all content produced by a model. This can be done after an image is generated using standard techniques [6] or can be baked into the synthesis by



watermarking all the training data [23]. Lastly, more care can be taken to ensure the provenance of training images by, for example, licensing images from trusted sources.

These approaches, of course, are not perfect: passive detection of AI-generated images is not perfect: a sophisticated adversary can remove a watermark, and provenance is not always available or completely reliable. Combined, however, these strategies will most likely mitigate some of the risk of data poisoning by significantly reducing the number of undesired images.

A few open questions remain. What about the underlying model or training data causes the data poisoning? Will data poisoning generalize across synthesis engines: Will, for example, Stable Diffusion retrained on DALL-E or Midjourney images exhibit the same type of model collapse? Can generative-AI systems be trained or modified to be resilient to this type of data poisoning? If it turns out to be difficult to prevent data poisoning, are there specific techniques or data sets that can accelerate healing?

## Acknowledgments

We thank Emily Cooper and Issa Alford for many helpful dinner-time discussions and inspiration.

## References

- [1] Sina Alemohammad, Josue Casco-Rodriguez, Lorenzo Luzi, Ahmed Imtiaz Humayun, Hossein Babaei, Daniel LeJeune, Ali Siahkoobi, and Richard G Baraniuk. Self-consuming generative models go mad. *arXiv:2307.01850*, 2023. 1
- [2] Nicholas Carlini, Matthew Jagielski, Christopher A Choquette-Choo, Daniel Paleka, Will Pearce, Hyrum Anderson, Andreas Terzis, Kurt Thomas, and Florian Tramèr. Poisoning web-scale training datasets is practical. *arXiv:2302.10149*, 2023. 8
- [3] Nicholas Carlini and Andreas Terzis. Poisoning and backdooring contrastive learning. *arXiv:2106.09667*, 2021. 8
- [4] Lucy Chai, David Bau, Ser-Nam Lim, and Phillip Isola. What makes fake images detectable? Understanding properties that generalize. In *European Conference on Computer Vision*, pages 103–120, 2020. 8
- [5] Riccardo Corvi, Davide Cozzolino, Giada Zingarini, Giovanni Poggi, Koki Nagano, and Luisa Verdoliva. On the detection of synthetic images generated by diffusion models. In *International Conference on Acoustics, Speech and Signal Processing*, pages 1–5. IEEE, 2023. 8
- [6] Ingemar J Cox, Matt L Miller, JMG Linnartz, and Ton Kalker. A review of watermarking principles and practices. *Digital Signal Processing for Multimedia Systems*, 2:461–482, 1999. 8

- [7] Diego Gragnaniello, Davide Cozzolino, Francesco Marra, Giovanni Poggi, and Luisa Verdoliva. Are GAN generated images easy to detect? A critical analysis of the state-of-the-art. In *IEEE International Conference on Multimedia and Expo*, pages 1–6, 2021. 8
- [8] Ryuichiro Hataya, Han Bao, and Hiromi Arai. Will large-scale generative models corrupt future datasets? In *IEEE/CVF International Conference on Computer Vision*, pages 20555–20565, 2023. 2, 4
- [9] Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. CLIPScore: A reference-free evaluation metric for image captioning. arXiv:2104.08718, 2021. 7
- [10] 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. 7
- [11] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4401–4410, 2019. 5
- [12] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of StyleGAN. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8110–8119, 2020. 1
- [13] Bo Liu, Fan Yang, Xiuli Bi, Bin Xiao, Weisheng Li, and Xinbo Gao. Detecting generated images by real images. In *European Conference on Computer Vision*, pages 95–110. Springer, 2022. 8
- [14] Gonzalo Martínez, Lauren Watson, Pedro Reviriego, José Alberto Hernández, Marc Juárez, and Rik Sarkar. Combining generative artificial intelligence (AI) and the internet: Heading towards evolution or degradation? arXiv:2303.01255, 2023. 1
- [15] Gonzalo Martínez, Lauren Watson, Pedro Reviriego, José Alberto Hernández, Marc Juárez, and Rik Sarkar. Towards understanding the interplay of generative artificial intelligence and the internet. arXiv:2306.06130, 2023. 1
- [16] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *International Conference on Machine Learning*, pages 8162–8171. PMLR, 2021. 2
- [17] 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. 7
- [18] Zohra Rezgui. Détection et classification de visages pour la description de l'égalité femme-homme dans les archives télévisuelles, 2019. 5

- [19] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10684–10695, 2022. 5
- [20] Sefik Ilkin Serengil and Alper Ozpinar. HyperExtended lightface: A facial attribute analysis framework. In *International Conference on Engineering and Emerging Technologies*, pages 1–4. IEEE, 2021. 5
- [21] Tamar Rott Shaham, Tali Dekel, and Tomer Michaeli. Singan: Learning a generative model from a single natural image. In *IEEE/CVF International Conference on Computer Vision*, pages 4570–4580, 2019. 8
- [22] Iliia Shumailov, Zakhar Shumaylov, Yiren Zhao, Yarin Gal, Nicolas Papernot, and Ross Anderson. The curse of recursion: Training on generated data makes models forget. arxiv:2305.17493, 2023. 1
- [23] Ning Yu, Vladislav Skripniuk, Sahar Abdelnabi, and Mario Fritz. Artificial fingerprinting for generative models: Rooting deepfake attribution in training data. In *IEEE/CVF International Conference on Computer Vision*, pages 14448–14457, 2021. 9
- [24] Xu Zhang, Svebor Karaman, and Shih-Fu Chang. Detecting and simulating artifacts in GAN fake images. In *IEEE International Workshop on Information Forensics and Security*, pages 1–6, 2019. 8