

# Some Challenges around Retraining Generative Models on Their Own Data

Quentin Bertrand

Inria – UJM – [QB3.github.io](https://QB3.github.io)



D. Ferbach



J. A. Bose



A. Duplessis



M. Jiralerspong



G. Gidel

# LAION-5B<sup>12</sup>

Backend url:

<https://knn5.laion.ai>

Index:

laion\_5B 

french cat



[Clip retrieval](#) works by converting the text query to a CLIP embedding, then using that embedding to query a knn index of clip image embeddings



french cat



french cat



How to tell if your feline is french. He wears a b...



イケメン猫モデル  
「トキ・ナンタケット」がかっこいい-  
NAVERまとめ



Hilarious pics of funny cats! [funnycatsgif.com](http://funnycatsgif.com)

Display captions

Display full

captions

Display similarities

Safe mode

Hide duplicate urls

Hide (near) duplicate images

Search over

[image](#) 



Hipster cat



網友挑戰「加幾筆畫出最創意貓咪圖片」，笑到岔氣之後我也手



cat in a suit Georgian sells tomatoes



French Bread Cat Loaf Metal Print



<sup>1</sup>C. Schuhmann et al. "Laion-5B: An open large-scale dataset for training next generation image-text models". In: *NeurIPS* (2022).

<sup>2</sup><https://paperswithcode.com/dataset/laion-5b>

# Generative Models Today

- ▶ Powerful deep generative models
  - ↪ e.g. Diffusion trained on LAION-5B
- ▶ Easy access (Midjourney, Stable Diffusion, DALL·E)

# Generative Models Today

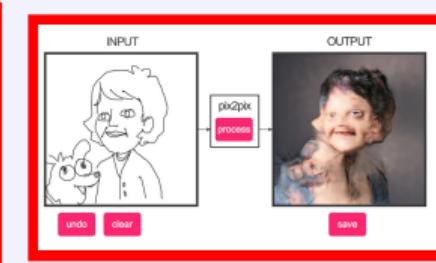
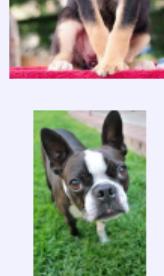
- ▶ Powerful deep generative models
  - ↪ e.g. Diffusion trained on LAION-5B
- ▶ Easy access (Midjourney, Stable Diffusion, DALL·E)
- ▶ Populates the WEB with **synthetically generated images**

# Generative Models Today

- ▶ Powerful deep generative models
  - ↪ e.g. Diffusion trained on LAION-5B
- ▶ Easy access (Midjourney, Stable Diffusion, DALL·E)
- ▶ Populates the WEB with **synthetically generated images**

# Inevitably Train on Synthetic Data

The LAION-5B<sup>3</sup> dataset already contains synthetically generated images<sup>4</sup>



<sup>3</sup>C. Schuhmann et al. "Laion-5B: An open large-scale dataset for training next generation image-text models". In: *NeurIPS* (2022).

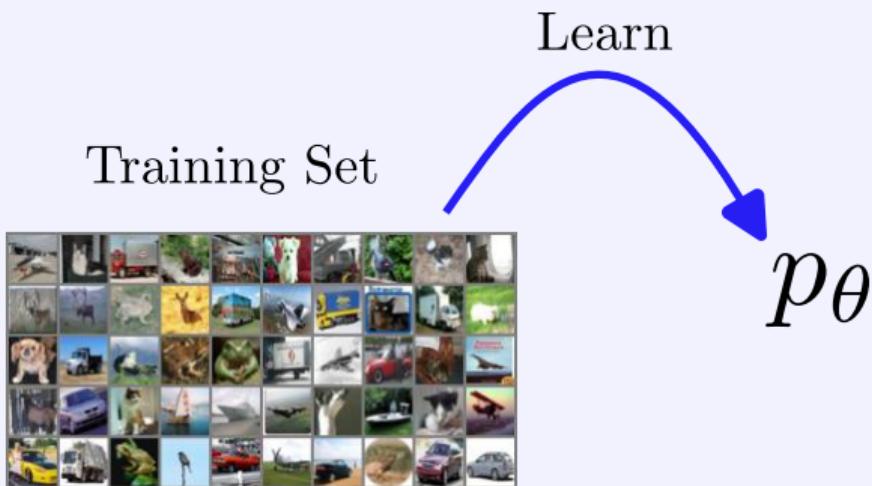
<sup>4</sup>S. Alemohammad et al. "Self-Consuming Generative Models Go MAD". In: *ICLR* (2024).

# Retraining Generative Models on their Own Data

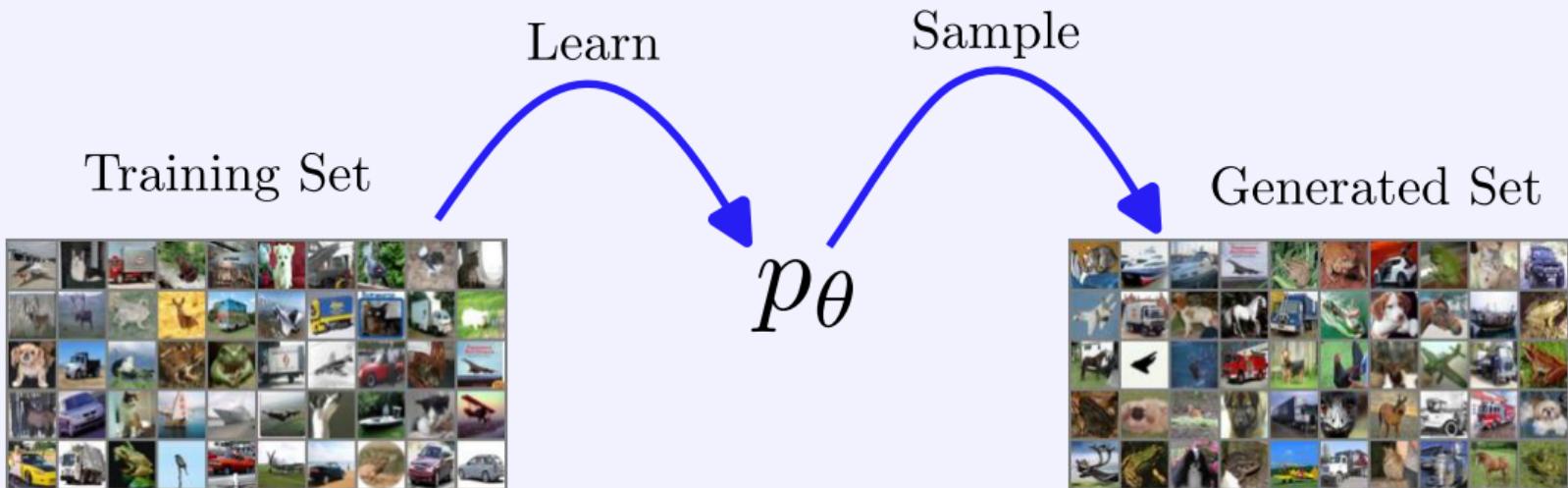
Training Set



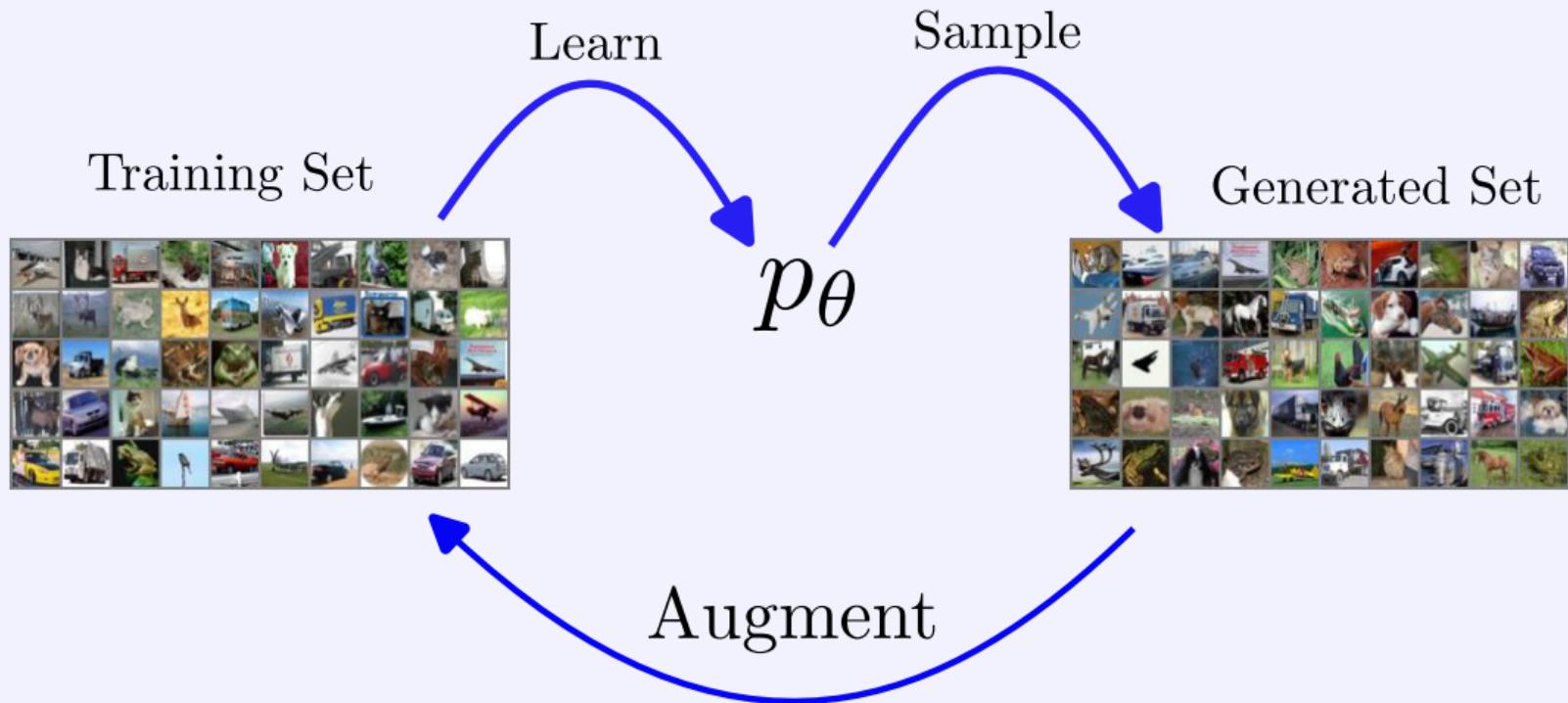
# Retraining Generative Models on their Own Data



# Retraining Generative Models on their Own Data



# Retraining Generative Models on their Own Data



# Training on Synthetic Data, Good or Bad?

## Iterative Retraining is Bad

- ▶ **The curse of recursion:** Training on generated data makes models forget<sup>a</sup>
- ▶ Self-Consuming Generative Models Go **MAD**<sup>b</sup>
- ▶ When **A.I. 's Output Is a Threat to A.I.** Itself (N.Y. Times article)

---

<sup>a</sup>I. Shumailov et al. "The Curse of Recursion: Training on Generated Data Makes Models Forget". In: (2023). arXiv: 2305.17493 [cs.LG].

<sup>b</sup>S. Alemohammad et al. "Self-Consuming Generative Models Go MAD". In: *ICLR* (2024).

# Training on Synthetic Data, Good or Bad?

## Iterative Retraining is Bad

- ▶ **The curse of recursion:** Training on generated data makes models forget<sup>a</sup>
- ▶ Self-Consuming Generative Models Go **MAD**<sup>b</sup>
- ▶ When **A.I. 's Output Is a Threat to A.I.** Itself (N.Y. Times article)

<sup>a</sup>I. Shumailov et al. "The Curse of Recursion: Training on Generated Data Makes Models Forget". In: (2023). arXiv: 2305.17493 [cs.LG].

<sup>b</sup>S. Alemdemehmed et al. "Self-Consuming Generative Models Go MAD". In: ICLR (2024).

Will generative models collapse?!

# Training on Synthetic Data, Good or Bad?

## Iterative Retraining is Bad

- ▶ **The curse of recursion:** Training on generated data makes models forget<sup>a</sup>
- ▶ Self-Consuming Generative Models Go **MAD**<sup>b</sup>
- ▶ When **A.I. 's Output Is a Threat to A.I.** Itself (N.Y. Times article)

<sup>a</sup>I. Shumailov et al. "The Curse of Recursion: Training on Generated Data Makes Models Forget". In: (2023). arXiv: 2305.17493 [cs.LG].

<sup>b</sup>S. Aleomohammad et al. "Self-Consuming Generative Models Go MAD". In: ICLR (2024).

## Will generative models collapse?!

## Training on Synthetic Data is Good

- ▶ Data augmentation for downstream tasks
  - ↪ Adversarial training<sup>a</sup>
  - ↪ Classification with imbalanced datasets<sup>b</sup>
  - ↪ Generative modelling<sup>c</sup>

<sup>a</sup>Z. Wang et al. "Better diffusion models further improve adversarial training". In: ICML. 2023.

<sup>b</sup>R. A. Hemmat et al. "Feedback-guided Data Synthesis for Imbalanced Classification". In: TMLR (2024).

<sup>c</sup>C. Gulcehre et al. "Reinforced self-training (REST) for language modeling". In: (2023). arXiv: 2308.08998 [cs.CL].

# Training on Synthetic Data, Good or Bad?

## Iterative Retraining is Bad

- ▶ **The curse of recursion:** Training on generated data makes models forget<sup>a</sup>
- ▶ Self-Consuming Generative Models Go **MAD**<sup>b</sup>
- ▶ When **A.I. 's Output Is a Threat to A.I.** Itself (N.Y. Times article)

<sup>a</sup>I. Shumailov et al. "The Curse of Recursion: Training on Generated Data Makes Models Forget". In: (2023). arXiv: 2305.17493 [cs.LG].

<sup>b</sup>S. Alemdoohammad et al. "Self-Consuming Generative Models Go MAD". In: ICLR (2024).

## Will generative models collapse?!

## Training on Synthetic Data is Good

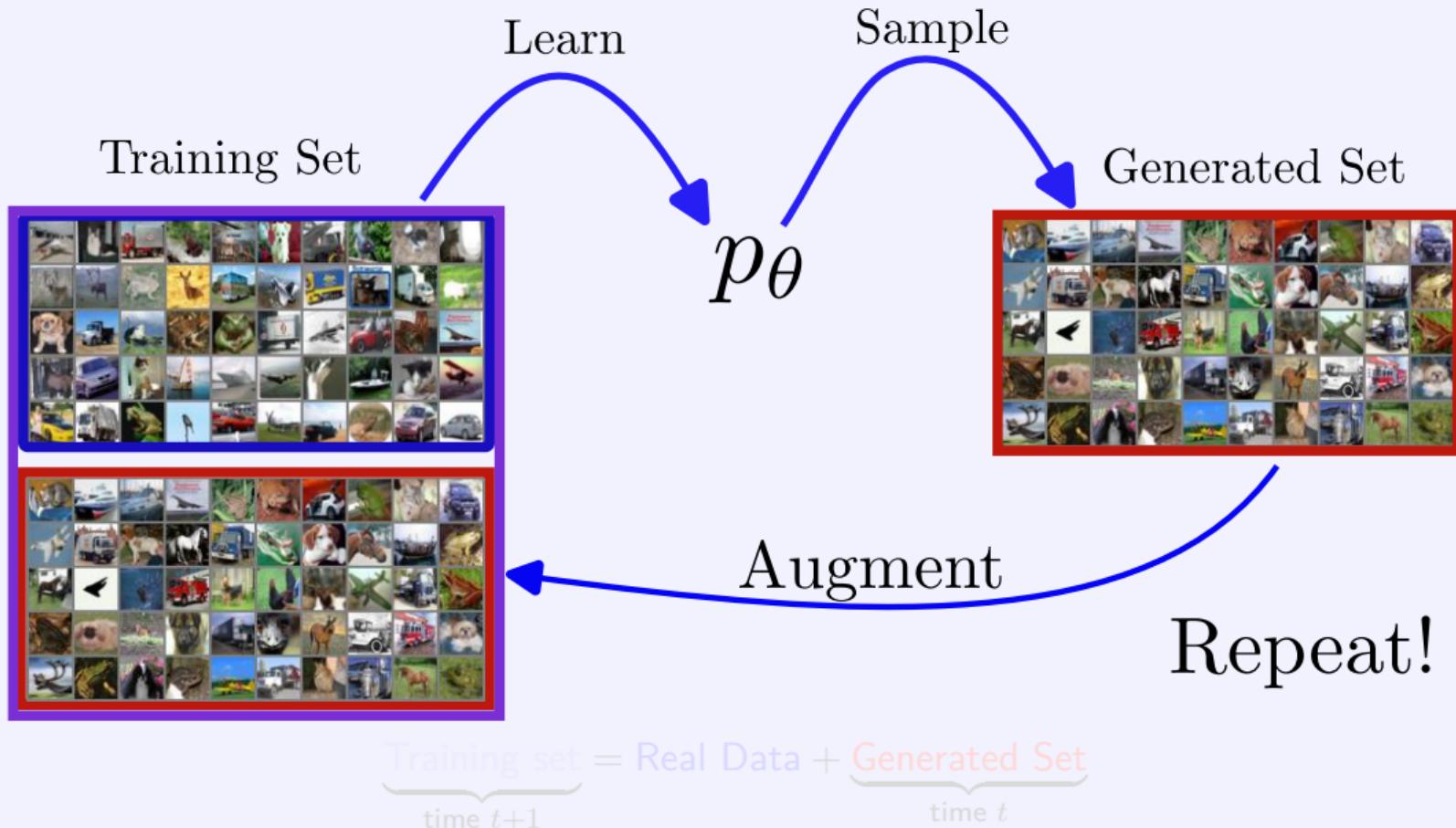
- ▶ Data augmentation for downstream tasks
  - ↪ Adversarial training<sup>a</sup>
  - ↪ Classification with imbalanced datasets<sup>b</sup>
  - ↪ Generative modelling<sup>c</sup>

<sup>a</sup>Z. Wang et al. "Better diffusion models further improve adversarial training". In: ICML. 2023.

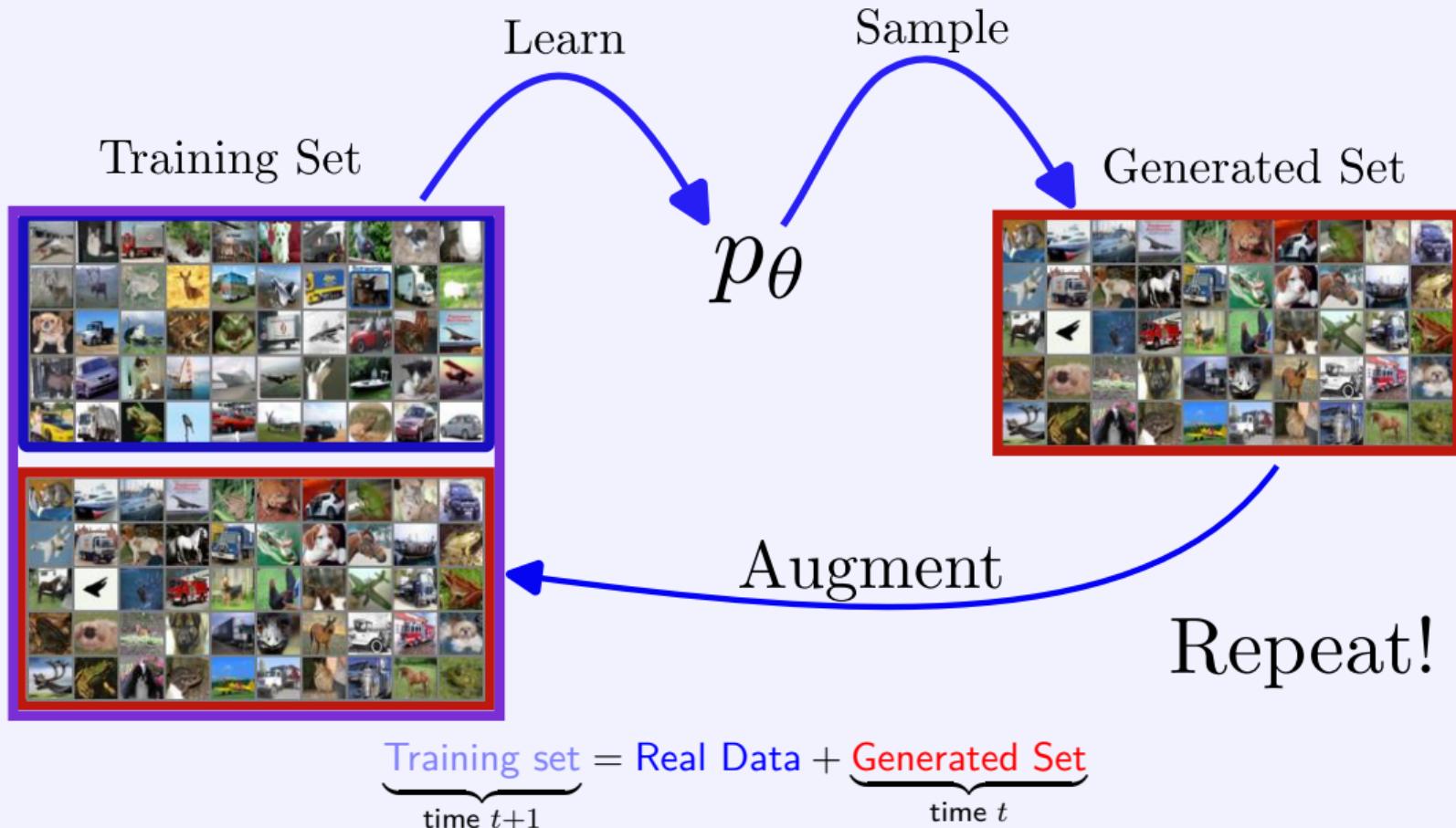
<sup>b</sup>R. A. Hemmat et al. "Feedback-guided Data Synthesis for Imbalanced Classification". In: TMLR (2024).

<sup>c</sup>C. Gulcehre et al. "Reinforced self-training (REST) for language modeling". In: (2023). arXiv: 2308.08998 [cs.CL].

# Iterative Retraining / Self-Consuming Generative Models



# Iterative Retraining / Self-Consuming Generative Models



# Setting

## Notation

- ▶  $\hat{p}_{\text{data}}$  Empirical data distribution
  - ↪  $n$  Data points
- ▶  $p$  Likelihood of the model
  - ↪ Parametrized by  $\theta^n \in \Theta$

## Iterative Retraining

$$p_0 \in \arg \max_{p \in \mathcal{P}_\Theta} \mathbb{E}_{x \sim \hat{p}_{\text{data}}} [\log p(x)]$$

# Setting

## Notation

- $\hat{p}_{\text{data}}$  Empirical data distribution  
     $\hookrightarrow n$  Data points
- $p$  Likelihood of the model  
     $\hookrightarrow$  Parametrized by  $\theta^n \in \Theta$

## Iterative Retraining

$$p_0 \in \arg \max_{p \in \mathcal{P}_\Theta} \mathbb{E}_{x \sim \hat{p}_{\text{data}}} [\log p(x)]$$

$$p_{t+1} \in \arg \max_{p \in \mathcal{P}_\Theta} \underbrace{\mathbb{E}_{x \sim \hat{p}_{\text{data}}} [\log p(x)]}_{\text{Real data}} + \lambda \cdot \underbrace{\mathbb{E}_{\tilde{x} \sim \hat{p}_t} [\log p(\tilde{x})]}_{\text{Synthetic data}}$$

# Setting

## Notation

- ▶  $\hat{p}_{\text{data}}$  Empirical data distribution  
↪  $n$  Data points
- ▶  $p$  Likelihood of the model  
↪ Parametrized by  $\theta^n \in \Theta$

## Iterative Retraining

$$p_0 \in \arg \max_{p \in \mathcal{P}_\Theta} \mathbb{E}_{x \sim \hat{p}_{\text{data}}} [\log p(x)]$$

$$p_{t+1} \in \arg \max_{p \in \mathcal{P}_\Theta} \underbrace{\mathbb{E}_{x \sim \hat{p}_{\text{data}}} \log p(x)}_{\text{Real data}} + \lambda \cdot \underbrace{\mathbb{E}_{\tilde{x} \sim \hat{p}_t} \log p(\tilde{x})}_{\text{Synthetic data}}$$

# Warm Up: Only Retrain on your Own Data 1/3

## Iterative Retraining

$$p_0^n \in \arg \max_{p \in \mathcal{P}_\Theta} \mathbb{E}_{x \sim \hat{p}_{\text{data}}} [\log p(x)]$$

$$p_{t+1}^n \in \arg \max_{p \in \mathcal{P}_\Theta} \underbrace{\mathbb{E}_{x \sim \hat{p}_{\text{data}}} \log p(x)}_{\text{Real data}} + \cancel{X} \cdot \underbrace{\mathbb{E}_{\tilde{x} \sim p_t} \log p(\tilde{x})}_{\text{Synthetic data}}$$

Q: What will happen?

# Warm Up: Only Retrain on your Own Data 1/3

## Iterative Retraining

$$p_0^n \in \arg \max_{p \in \mathcal{P}_\Theta} \mathbb{E}_{x \sim \hat{p}_{\text{data}}} [\log p(x)]$$

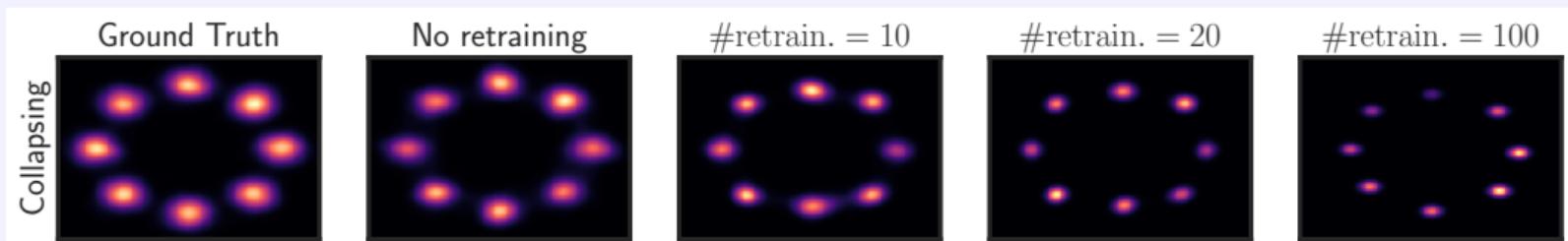
$$p_{t+1}^n \in \arg \max_{p \in \mathcal{P}_\Theta} \underbrace{\mathbb{E}_{x \sim \hat{p}_{\text{data}}} \log p(x)}_{\text{Real data}} + X \cdot \underbrace{\mathbb{E}_{\tilde{x} \sim p_t} \log p(\tilde{x})}_{\text{Synthetic data}}$$

Q: What will happen?

## Warm Up: Only Retrain on your Own Data 2/3

Q: What will happen?

A: Mode Collapse



### Setup

- ▶ Data: 8 Gaussians,  $x \in \mathbb{R}^2$
- ▶ Algorithm: Diffusion (DDPM<sup>a</sup>)

<sup>a</sup>J. Ho, A. Jain, and P. Abbeel. "Denoising diffusion probabilistic models". In: *NeurIPS* (2020).

## Warm Up: Only Retrain on your Own Data 3/3

Single unidimensional Gaussian, unbiased estimator

Data:  $\textcolor{blue}{x}_j^0 = \mu_0 + \sigma_0 Z_j$ , with  $Z_j \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}_{0,1}$ ,  $1 \leq j \leq n$

Learning step: 
$$\begin{cases} \mu_{t+1} &= \frac{1}{n} \sum_j \tilde{\mathbf{x}}_j^t \\ \sigma_{t+1}^2 &= \frac{1}{n-1} \sum_j (\tilde{\mathbf{x}}_j^t - \mu_{t+1})^2 \end{cases}$$

Sampling step:  $\tilde{\mathbf{x}}_j^{t+1} = \mu_{t+1} + \sigma_{t+1} \cdot Z_j^{t+1}$ , with  $Z_j^{t+1} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}_{0,1}$ ,  $1 \leq j \leq n$

Result

$$\mathbb{E}(\sigma_t) \leq \alpha^t \sigma_0 \underset{t \rightarrow +\infty}{\longrightarrow} 0, \quad 0 \leq \alpha < 1$$

## Warm Up: Only Retrain on your Own Data 3/3

Single unidimensional Gaussian, unbiased estimator

Data:  $\textcolor{blue}{x}_j^0 = \mu_0 + \sigma_0 Z_j$ , with  $Z_j \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}_{0,1}$ ,  $1 \leq j \leq n$

Learning step: 
$$\begin{cases} \mu_{t+1} &= \frac{1}{n} \sum_j \tilde{\mathbf{x}}_j^t \\ \sigma_{t+1}^2 &= \frac{1}{n-1} \sum_j (\tilde{\mathbf{x}}_j^t - \mu_{t+1})^2 \end{cases}$$

Sampling step:  $\tilde{\mathbf{x}}_j^{t+1} = \mu_{t+1} + \sigma_{t+1} \cdot Z_j^{t+1}$ , with  $Z_j^{t+1} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}_{0,1}$ ,  $1 \leq j \leq n$

Result

$$\mathbb{E}(\sigma_t) \leq \alpha^t \sigma_0 \underset{t \rightarrow +\infty}{\longrightarrow} 0, \quad 0 \leq \alpha < 1$$

Same type of results holds for a single multidimensional Gaussian

## Warm Up: Only Retrain on your Own Data 3/3

Single unidimensional Gaussian, unbiased estimator

Data:  $\textcolor{blue}{x}_j^0 = \mu_0 + \sigma_0 Z_j$ , with  $Z_j \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}_{0,1}$ ,  $1 \leq j \leq n$

Learning step: 
$$\begin{cases} \mu_{t+1} &= \frac{1}{n} \sum_j \tilde{\mathbf{x}}_j^t \\ \sigma_{t+1}^2 &= \frac{1}{n-1} \sum_j (\tilde{\mathbf{x}}_j^t - \mu_{t+1})^2 \end{cases}$$

Sampling step:  $\tilde{\mathbf{x}}_j^{t+1} = \mu_{t+1} + \sigma_{t+1} \cdot Z_j^{t+1}$ , with  $Z_j^{t+1} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}_{0,1}$ ,  $1 \leq j \leq n$

Result

$$\mathbb{E}(\sigma_t) \leq \alpha^t \sigma_0 \underset{t \rightarrow +\infty}{\longrightarrow} 0, \quad 0 \leq \alpha < 1$$

Same type of results holds for a single multidimensional Gaussian

# General Case

## Iterative Retraining

$$p_0 \in \arg \max_{p \in \mathcal{P}_\Theta} \mathbb{E}_{x \sim \hat{p}_{\text{data}}} [\log p(x)]$$

$$p_{t+1} \in \arg \max_{p \in \mathcal{P}_\Theta} \underbrace{\mathbb{E}_{x \sim \hat{p}_{\text{data}}} [\log p(\textcolor{blue}{x})]}_{\text{Real data}} + \underbrace{\lambda \cdot \mathbb{E}_{\tilde{x} \sim p_t} [\log p(\tilde{\textcolor{red}{x}})}_{\text{Synthetic data}} := \mathcal{G}(p_t)$$

Idea

- ▶ Fixed-point iteration  $p_{t+1} = \mathcal{G}(p_t)$

# General Case

## Iterative Retraining

$$p_0 \in \arg \max_{p \in \mathcal{P}_\Theta} \mathbb{E}_{x \sim \hat{p}_{\text{data}}} [\log p(x)]$$

$$\boldsymbol{p}_{t+1} \in \arg \max_{p \in \mathcal{P}_\Theta} \underbrace{\mathbb{E}_{\boldsymbol{x} \sim \hat{p}_{\text{data}}} [\log p(\boldsymbol{x})]}_{\text{Real data}} + \lambda \cdot \underbrace{\mathbb{E}_{\tilde{\boldsymbol{x}} \sim \boldsymbol{p}_t} [\log p(\tilde{\boldsymbol{x}})]}_{\text{Synthetic data}} := \mathcal{G}(\boldsymbol{p}_t)$$

## Idea

- ▶ Fixed-point iteration  $\boldsymbol{p}_{t+1} = \mathcal{G}(\boldsymbol{p}_t)$
- ▶  $p_{\text{data}}$  is a fixed point of  $\mathcal{G}$ 
  - ▶ Study the stability of  $\mathcal{G}$  around  $p_{\text{data}}$

# General Case

## Iterative Retraining

$$p_0 \in \arg \max_{p \in \mathcal{P}_\Theta} \mathbb{E}_{x \sim \hat{p}_{\text{data}}} [\log p(x)]$$

$$p_{t+1} \in \arg \max_{p \in \mathcal{P}_\Theta} \underbrace{\mathbb{E}_{x \sim \hat{p}_{\text{data}}} [\log p(\textcolor{blue}{x})]}_{\text{Real data}} + \lambda \cdot \underbrace{\mathbb{E}_{\tilde{x} \sim p_t} [\log p(\tilde{\textcolor{red}{x}})}_{\text{Synthetic data}} := \mathcal{G}(p_t)$$

## Idea

- ▶ Fixed-point iteration  $p_{t+1} = \mathcal{G}(p_t)$
- ▶  $\hat{p}_{\text{data}}$  is a fixed point of  $\mathcal{G}$ 
  - ▶ Study the stability of  $\mathcal{G}$  around  $\hat{p}_{\text{data}}$
- ▶ Link with performative prediction!

# General Case

## Iterative Retraining

$$p_0 \in \arg \max_{p \in \mathcal{P}_\Theta} \mathbb{E}_{x \sim \hat{p}_{\text{data}}} [\log p(x)]$$

$$\boldsymbol{p}_{t+1} \in \arg \max_{p \in \mathcal{P}_\Theta} \underbrace{\mathbb{E}_{x \sim \hat{p}_{\text{data}}} \log p(\boldsymbol{x})}_{\text{Real data}} + \lambda \cdot \underbrace{\mathbb{E}_{\tilde{\boldsymbol{x}} \sim \boldsymbol{p}_t} \log p(\tilde{\boldsymbol{x}})}_{\text{Synthetic data}} := \mathcal{G}(\boldsymbol{p}_t)$$

## Idea

- ▶ Fixed-point iteration  $\boldsymbol{p}_{t+1} = \mathcal{G}(\boldsymbol{p}_t)$
- ▶  $\hat{p}_{\text{data}}$  is a fixed point of  $\mathcal{G}$ 
  - ▶ Study the stability of  $\mathcal{G}$  around  $\hat{p}_{\text{data}}$
- ▶ Link with performative prediction!

## Retrain of Generative Models: Informal

### Assumptions

- ▶ Regularity of the log-likelihood
  - ↪ Local Lipschitzness and strong convexity
- ▶ The first generative model is "good enough"
  - ↪  $\mathcal{W}(p_{\text{data}}, p_{\theta_0}) < \epsilon$

## Retrain of Generative Models: Informal

### Assumptions

- ▶ Regularity of the log-likelihood
  - ↪ Local Lipschitzness and strong convexity
- ▶ The first generative model is "good enough"
  - ↪  $\mathcal{W}(p_{\text{data}}, p_{\theta_0}) < \epsilon$
- ▶ Infinite Data

# Retrain of Generative Models: Informal

## Assumptions

- ▶ Regularity of the log-likelihood
  - ↪ Local Lipschitzness and strong convexity
- ▶ The first generative model is "good enough"
  - ↪  $\mathcal{W}(p_{\text{data}}, p_{\theta_0}) < \epsilon$
- ▶ Infinite Data

## Results

- ▶ Regularity + good enough model + infinite data

---

<sup>10</sup> Q. Beinard et al. "On the stability of iterative retraining of generative models on their own data". In: ICML 2021.

<sup>11</sup> D. Ferbach et al. "Self-Consuming Generative Models with Curated Data Provably Optimize Human Preferences". In: NeurIPS 2021.

# Retrain of Generative Models: Informal

## Assumptions

- ▶ Regularity of the log-likelihood
  - ↪ Local Lipschitzness and strong convexity
- ▶ The first generative model is "good enough"
  - ↪  $\mathcal{W}(p_{\text{data}}, p_{\theta_0}) < \epsilon$
- ▶ Infinite Data

## Results

- ▶ Regularity + good enough model + infinite data
  - ⇒ Stability: if  $\theta_0$  is close enough to  $\theta^*$ , then  $\text{KL}(p_{\theta^*} \| p_{\theta_t}) \rightarrow 0$  linearly<sup>a b</sup>

---

<sup>a</sup>Q. Bertrand et al. "On the stability of iterative retraining of generative models on their own data". In: *ICLR* (2024).

<sup>b</sup>D. Ferbach et al. "Self-Consuming Generative Models with Curated Data Provably Optimize Human Preferences". In: *NeurIPS* (2024).

# Retrain of Generative Models: Informal

## Assumptions

- ▶ Regularity of the log-likelihood
  - ↪ Local Lipschitzness and strong convexity
- ▶ The first generative model is "good enough"
  - ↪  $\mathcal{W}(p_{\text{data}}, p_{\theta_0}) < \epsilon$
- ▶ Infinite Data

## Results

- ▶ Regularity + good enough model + infinite data
  - ⇒ Stability: if  $\theta_0$  is close enough to  $\theta^*$ , then  $\text{KL}(p_{\theta^*} \| p_{\theta_t}) \rightarrow 0$  linearly<sup>a</sup><sup>b</sup>
  - ↪ Interplay between prop. of real data and  $\mathcal{W}(p_{\text{data}}, p_{\theta_0}) < \epsilon$

---

<sup>a</sup>Q. Bertrand et al. "On the stability of iterative retraining of generative models on their own data". In: *ICLR* (2024).

<sup>b</sup>D. Ferbach et al. "Self-Consuming Generative Models with Curated Data Provably Optimize Human Preferences". In: *NeurIPS* (2024).

# Retrain of Generative Models: Informal

## Assumptions

- ▶ Regularity of the log-likelihood
  - ↪ Local Lipschitzness and strong convexity
- ▶ The first generative model is "good enough"
  - ↪  $\mathcal{W}(p_{\text{data}}, p_{\theta_0}) < \epsilon$
- ▶ Infinite Data

## Results

- ▶ Regularity + good enough model + infinite data
  - ⇒ Stability: if  $\theta_0$  is close enough to  $\theta^*$ , then  $\text{KL}(p_{\theta^*} \| p_{\theta_t}) \rightarrow 0$  linearly<sup>a</sup><sup>b</sup>
  - ↪ Interplay between prop. of real data and  $\mathcal{W}(p_{\text{data}}, p_{\theta_0}) < \epsilon$
  - ↪ Finite sample extension

<sup>a</sup>Q. Bertrand et al. "On the stability of iterative retraining of generative models on their own data". In: *ICLR* (2024).

<sup>b</sup>D. Ferbach et al. "Self-Consuming Generative Models with Curated Data Provably Optimize Human Preferences". In: *NeurIPS* (2024).

# Retrain of Generative Models: Informal

## Assumptions

- ▶ Regularity of the log-likelihood
  - ↪ Local Lipschitzness and strong convexity
- ▶ The first generative model is "good enough"
  - ↪  $\mathcal{W}(p_{\text{data}}, p_{\theta_0}) < \epsilon$
- ▶ Infinite Data

## Results

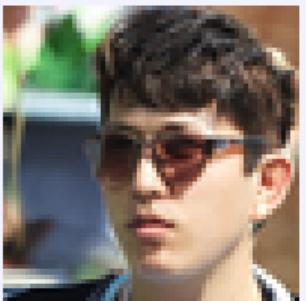
- ▶ Regularity + good enough model + infinite data
  - ⇒ Stability: if  $\theta_0$  is close enough to  $\theta^*$ , then  $\text{KL}(p_{\theta^*} \| p_{\theta_t}) \rightarrow 0$  linearly<sup>a</sup><sup>b</sup>
  - ↪ Interplay between prop. of real data and  $\mathcal{W}(p_{\text{data}}, p_{\theta_0}) < \epsilon$
  - ↪ Finite sample extension

---

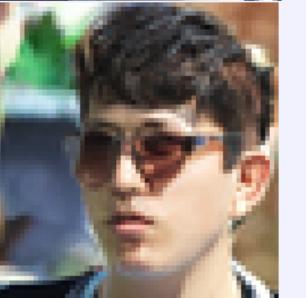
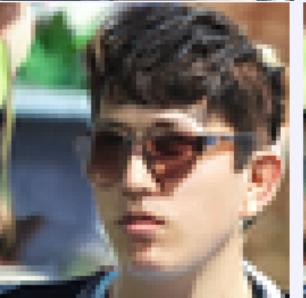
<sup>a</sup>Q. Bertrand et al. "On the stability of iterative retraining of generative models on their own data". In: *ICLR* (2024).

<sup>b</sup>D. Ferbach et al. "Self-Consuming Generative Models with Curated Data Provably Optimize Human Preferences". In: *NeurIPS* (2024).

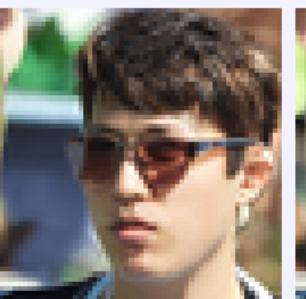
Fully synth.



$\lambda = 0.5$



$\lambda = 0$



0 retrain.

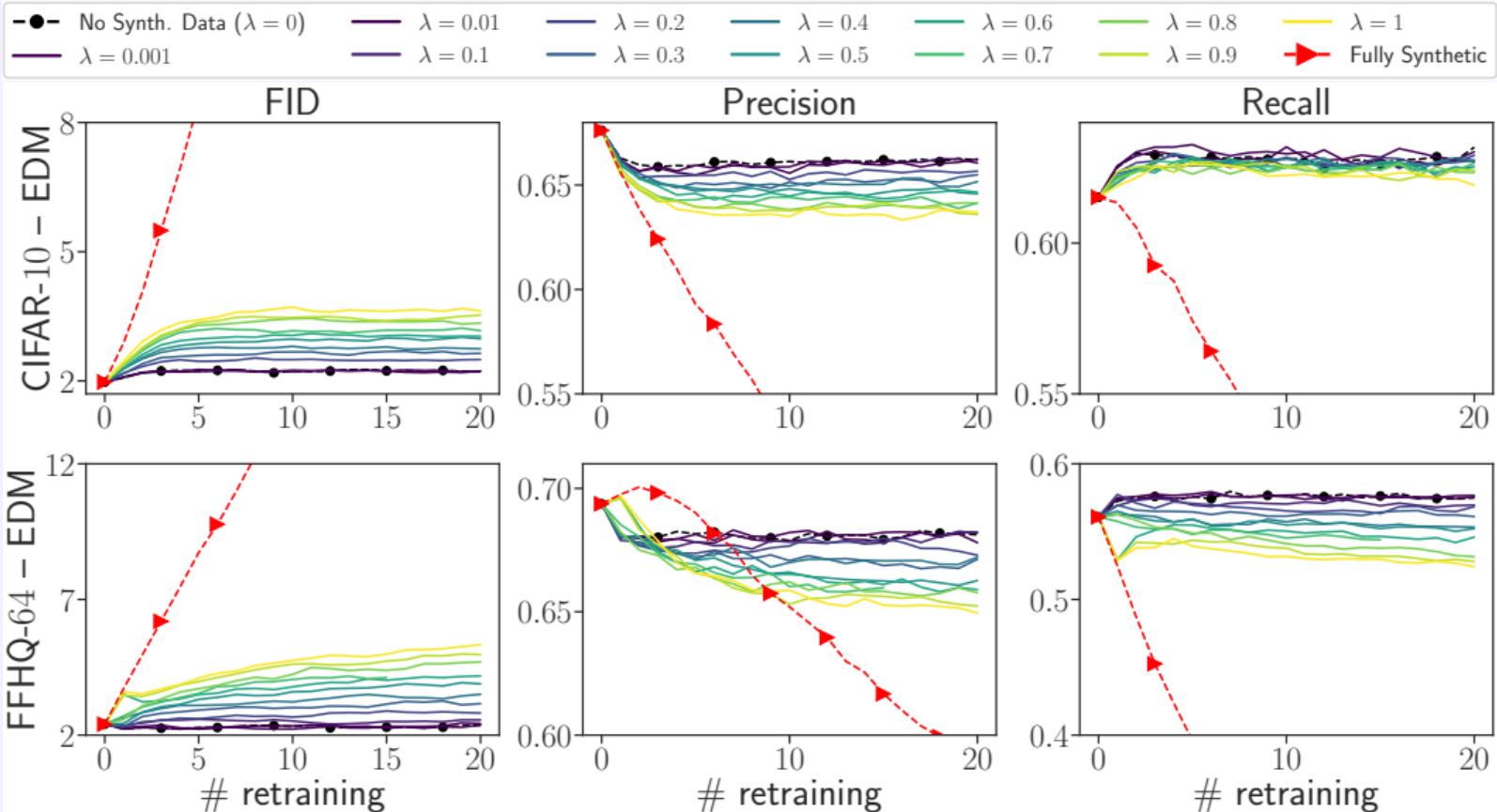
5 retrain.

10 retrain.

15 retrain.

20 retrain.

# Experiments



# Intermediate Conclusions

## Self-consuming generative models

- ▶ No collapse/MADness (if "enough" real data)<sup>a b</sup>

---

<sup>a</sup>Q. Bertrand et al. "On the stability of iterative retraining of generative models on their own data". In: *ICLR* (2024).

<sup>b</sup>R. Hataya, H. Bao, and H.J. Arai. "Will Large-scale Generative Models Corrupt Future Datasets?" In: *ICCV*. 2023.

# Intermediate Conclusions

## Self-consuming generative models

- ▶ No collapse/MADness (if "enough" real data)<sup>a b</sup>
- ▶ No improvements either

---

<sup>a</sup>Q. Bertrand et al. "On the stability of iterative retraining of generative models on their own data". In: *ICLR* (2024).

<sup>b</sup>R. Hataya, H. Bao, and HJ. Arai. "Will Large-scale Generative Models Corrupt Future Datasets?" In: *ICCV*. 2023.

# Intermediate Conclusions

## Self-consuming generative models

- ▶ No collapse/MADness (if "enough" real data)<sup>a b</sup>
- ▶ No improvements either

---

<sup>a</sup>Q. Bertrand et al. "On the stability of iterative retraining of generative models on their own data". In: *ICLR* (2024).

<sup>b</sup>R. Hataya, H. Bao, and HJ. Arai. "Will Large-scale Generative Models Corrupt Future Datasets?" In: *ICCV*. 2023.

What if you already retrain on curated/filtered synthetic data?

# Intermediate Conclusions

## Self-consuming generative models

- ▶ No collapse/MADness (if "enough" real data)<sup>a b</sup>
- ▶ No improvements either

---

<sup>a</sup>Q. Bertrand et al. "On the stability of iterative retraining of generative models on their own data". In: *ICLR* (2024).

<sup>b</sup>R. Hataya, H. Bao, and HJ. Arai. "Will Large-scale Generative Models Corrupt Future Datasets?" In: *ICCV*. 2023.

What if you already retrain on curated/filtered synthetic data?

# Already Training on Curated/Filtered Data 1/2

## LAION-Aesthetics<sup>a</sup>

<sup>a</sup><https://github.com/LAION-AI/laion-datasets/blob/main/laion-aesthetic.md>

- ▶ Filter the LAION-5B dataset<sup>a</sup>
  - ↪ Filter: reward model<sup>b</sup>

---

<sup>a</sup>C. Schuhmann et al. "Laion-5B: An open large-scale dataset for training next generation image-text models". In: *NeurIPS* (2022).

<sup>b</sup><https://github.com/LAION-AI/aesthetic-predictor>

<sup>c</sup>J. D. Pressman, K. Crowson, and Simulacra Captions Contributors. *Simulacra Aesthetic Captions*. Version 1.0. url  
<https://github.com/JD-P/simulacra-aesthetic-captions>. Stability AI, 2022.

# Already Training on Curated/Filtered Data 1/2

## LAION-Aesthetics<sup>a</sup>

<sup>a</sup><https://github.com/LAION-AI/laion-datasets/blob/main/laion-aesthetic.md>

► Filter the LAION-5B dataset<sup>a</sup>

    ↳ Filter: reward model<sup>b</sup>

        ↳ Trained on synthetic data "Simulacra Aesthetic Captions" dataset<sup>c</sup>

---

<sup>a</sup>C. Schuhmann et al. "Laion-5B: An open large-scale dataset for training next generation image-text models". In: *NeurIPS* (2022).

<sup>b</sup><https://github.com/LAION-AI/aesthetic-predictor>

<sup>c</sup>J. D. Pressman, K. Crowson, and Simulacra Captions Contributors. *Simulacra Aesthetic Captions*. Version 1.0. url <https://github.com/JD-P/simulacra-aesthetic-captions>. Stability AI, 2022.

# Already Training on Curated/Filtered Data 1/2

## LAION-Aesthetics<sup>a</sup>

<sup>a</sup><https://github.com/LAION-AI/laion-datasets/blob/main/laion-aesthetic.md>

- ▶ Filter the LAION-5B dataset<sup>a</sup>
  - ↪ Filter: reward model<sup>b</sup>
    - ↪ Trained on synthetic data "Simulacra Aesthetic Captions" dataset<sup>c</sup>
- ▶ Filter 8M/120M sample subset of LAION-5B

---

<sup>a</sup>C. Schuhmann et al. "Laion-5B: An open large-scale dataset for training next generation image-text models". In: *NeurIPS* (2022).

<sup>b</sup><https://github.com/LAION-AI/aesthetic-predictor>

<sup>c</sup>J. D. Pressman, K. Crowson, and Simulacra Captions Contributors. *Simulacra Aesthetic Captions*. Version 1.0. url <https://github.com/JD-P/simulacra-aesthetic-captions>. Stability AI, 2022.

# Already Training on Curated/Filtered Data 1/2

## LAION-Aesthetics<sup>a</sup>

<sup>a</sup><https://github.com/LAION-AI/laion-datasets/blob/main/laion-aesthetic.md>

- ▶ Filter the LAION-5B dataset<sup>a</sup>
  - ↪ Filter: reward model<sup>b</sup>
    - ↪ Trained on synthetic data "Simulacra Aesthetic Captions" dataset<sup>c</sup>
- ▶ Filter 8M/120M sample subset of LAION-5B

<sup>a</sup>C. Schuhmann et al. "Laion-5B: An open large-scale dataset for training next generation image-text models". In: *NeurIPS* (2022).

<sup>b</sup><https://github.com/LAION-AI/aesthetic-predictor>

<sup>c</sup>J. D. Pressman, K. Crowson, and Simulacra Captions Contributors. *Simulacra Aesthetic Captions*. Version 1.0. url <https://github.com/JD-P/simulacra-aesthetic-captions>. Stability AI, 2022.

# Already Training on Curated/Filtered Data 2/2

## JourneyDB Dataset<sup>a</sup>

---

<sup>a</sup>P. Junting et al. "JourneyDB: A Benchmark for Generative Image Understanding". In: *NeurIPS* (2023).

- ▶ Midjourney<sup>a</sup>: text-to-image model
  - ▶ Midjourney "*discord*" server
- 

<sup>a</sup>Midjourney. <https://www.midjourney.com/home/>. Version 5.2. Accessed: 2023-09-09. 2023.

# Already Training on Curated/Filtered Data 2/2

## JourneyDB Dataset<sup>a</sup>

---

<sup>a</sup>P. Junting et al. "JourneyDB: A Benchmark for Generative Image Understanding". In: *NeurIPS* (2023).

- ▶ Midjourney<sup>a</sup>: text-to-image model
- ▶ Midjourney "*discord*" server
  - ↪ User can request prompts

---

<sup>a</sup>Midjourney. <https://www.midjourney.com/home/>. Version 5.2. Accessed: 2023-09-09. 2023.

# Already Training on Curated/Filtered Data 2/2

## JourneyDB Dataset<sup>a</sup>

---

<sup>a</sup>P. Junting et al. "JourneyDB: A Benchmark for Generative Image Understanding". In: *NeurIPS* (2023).

- ▶ Midjourney<sup>a</sup>: text-to-image model
  - ▶ Midjourney "discord" server
    - ↪ User can request prompts
    - ↪ Midjourney proposes  $K = 4$  images
- 

<sup>a</sup>Midjourney. <https://www.midjourney.com/home/>. Version 5.2. Accessed: 2023-09-09. 2023.

# Already Training on Curated/Filtered Data 2/2

## JourneyDB Dataset<sup>a</sup>

---

<sup>a</sup>P. Junting et al. "JourneyDB: A Benchmark for Generative Image Understanding". In: *NeurIPS* (2023).

- ▶ Midjourney<sup>a</sup>: text-to-image model
  - ▶ Midjourney "discord" server
    - ↪ User can request prompts
    - ↪ Midjourney proposes  $K = 4$  images
    - ↪ Users vote for which image to upscale
- 

<sup>a</sup>Midjourney. <https://www.midjourney.com/home/>. Version 5.2. Accessed: 2023-09-09. 2023.

# Already Training on Curated/Filtered Data 2/2

## JourneyDB Dataset<sup>a</sup>

<sup>a</sup>P. Junting et al. "JourneyDB: A Benchmark for Generative Image Understanding". In: *NeurIPS* (2023).

- ▶ Midjourney<sup>a</sup>: text-to-image model
- ▶ Midjourney "discord" server
  - ↪ User can request prompts
  - ↪ Midjourney proposes  $K = 4$  images
  - ↪ Users vote for which image to upscale

<sup>a</sup>Midjourney. <https://www.midjourney.com/home/>. Version 5.2. Accessed: 2023-09-09. 2023.

## Remarks

- ▶ No access to pairwise comparisons
  - ↪ Only access to the "winning" samples
  - ↪ As opposed to RHLF<sup>ab</sup>

<sup>a</sup>D. M. Ziegler et al. "Fine-tuning language models from human preferences". In: *arXiv preprint arXiv:1909.08593* (2019).

<sup>b</sup>R. Rafailov et al. "Direct preference optimization: Your language model is secretly a reward model". In: *NeurIPS* (2024).

# Already Training on Curated/Filtered Data 2/2

## JourneyDB Dataset<sup>a</sup>

<sup>a</sup>P. Junting et al. "JourneyDB: A Benchmark for Generative Image Understanding". In: *NeurIPS* (2023).

- ▶ Midjourney<sup>a</sup>: text-to-image model
- ▶ Midjourney "discord" server
  - ↪ User can request prompts
  - ↪ Midjourney proposes  $K = 4$  images
  - ↪ Users vote for which image to upscale

<sup>a</sup>Midjourney. <https://www.midjourney.com/home/>. Version 5.2. Accessed: 2023-09-09. 2023.

## Remarks

- ▶ No access to pairwise comparisons
  - ↪ Only access to the "winning" samples
  - ↪ As opposed to RHLF<sup>ab</sup>

<sup>a</sup>D. M. Ziegler et al. "Fine-tuning language models from human preferences". In: *arXiv preprint arXiv:1909.08593* (2019).

<sup>b</sup>R. Rafailov et al. "Direct preference optimization: Your language model is secretly a reward model". In: *NeurIPS* (2024).

# A Simple Curation Model 1/2

1. User picks prompts  $y \sim p_{user}(y)$
  2. Sample:  $x_1, x_2, x_3, x_4 \sim p_t(x | y)$
- REPEAT
5. Train the model  $p_{t+1}$  on the dataset  $\mathcal{D}_t$
  4. Only the **upscaled** samples are in the dataset  $\mathcal{D}_t$



3. Sample  $x_k$  is upscaled by User with probability:

$$\frac{e^{r(x_k)}}{\sum_{i=1}^4 e^{r(x_i)}}$$



# A Simple Curation Model 2/2

## Curation Model

- ▶ Suppose the existence of a reward model  $r$ 
  - ↪ score  $r(x)$  to each sample  $x$
- ▶
  - Sample  $\tilde{\mathbf{x}}_1 \sim p_t, \dots, \tilde{\mathbf{x}}_K \sim p_t$ , i.i.d.
- ▶
  - Pick  $\hat{\mathbf{x}} \sim \mathcal{BT}(\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_K)$ , i.e.,
$$\mathbb{P}(\hat{\mathbf{x}} = \tilde{\mathbf{x}}_k | \tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_K) = \frac{e^{r(\tilde{\mathbf{x}}_k)}}{\sum_{j=1}^K e^{r(\tilde{\mathbf{x}}_j)}}, 1 \leq k \leq K$$

## Iterative Retraining

$$p_{t+1} = \arg \max_{p \in \mathcal{P}} \mathbb{E}_{x \sim \hat{p}_{\text{data}}} [\log p(x)] + \lambda \cdot \mathbb{E}_{\substack{\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_K \sim p_t \\ \hat{\mathbf{x}} \sim \mathcal{BT}(\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_K)}} [\log p(\hat{\mathbf{x}})]$$

# A Simple Curation Model 2/2

## Curation Model

- ▶ Suppose the existence of a reward model  $r$ 
  - ↪ score  $r(x)$  to each sample  $x$
- ▶
  - Sample  $\tilde{\mathbf{x}}_1 \sim p_t, \dots, \tilde{\mathbf{x}}_K \sim p_t$ , i.i.d.
  - ▶ Pick  $\hat{\mathbf{x}} \sim \mathcal{BT}(\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_K)$ , i.e.,
$$\mathbb{P}(\hat{\mathbf{x}} = \tilde{\mathbf{x}}_k | \tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_K) = \frac{e^{r(\tilde{\mathbf{x}}_k)}}{\sum_{j=1}^K e^{r(\tilde{\mathbf{x}}_j)}}, 1 \leq k \leq K$$

## Iterative Retraining

$$p_{t+1} = \arg \max_{p \in \mathcal{P}} \mathbb{E}_{x \sim \hat{p}_{\text{data}}} [\log p(x)] + \lambda \cdot \mathbb{E}_{\substack{\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_K \sim p_t \\ \hat{\mathbf{x}} \sim \mathcal{BT}(\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_K)}} [\log p(\hat{\mathbf{x}})]$$

# Self-Consuming Loop

$$p_{t+1} = \arg \max_{p \in \mathcal{P}} \mathbb{E}_{x \sim \hat{p}_{\text{data}}} [\log p(x)] + \lambda \cdot \mathbb{E}_{\hat{\mathbf{x}}_1, \dots, \hat{\mathbf{x}}_K \sim p_t} [\log p(\hat{\mathbf{x}})]$$

Q: What will happen?

# Self-Consuming Loop

$$p_{t+1} = \arg \max_{p \in \mathcal{P}} \mathbb{E}_{x \sim \hat{p}_{\text{data}}} [\log p(x)] + \lambda \cdot \mathbb{E}_{\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_K \sim p_t} [\log p(\hat{\mathbf{x}})]$$

Q: What will happen?

## Results

$$\mathbb{E}_{p_t} [e^{r(x)}] \xrightarrow{t \rightarrow \infty} e^{r_*} \quad \text{and} \quad \text{Var}_{p_t} [e^{r(x)}] \xrightarrow{t \rightarrow \infty} 0 .$$

# Self-Consuming Loop

$$p_{t+1} = \arg \max_{p \in \mathcal{P}} \mathbb{E}_{x \sim \hat{p}_{\text{data}}} [\log p(x)] + \lambda \cdot \mathbb{E}_{\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_K \sim p_t} [\log p(\hat{\mathbf{x}})]$$

Q: What will happen?

## Results

$$\mathbb{E}_{p_t} [e^{r(x)}] \xrightarrow{t \rightarrow \infty} e^{r^*} \quad \text{and} \quad \text{Var}_{p_t} [e^{r(x)}] \xrightarrow{t \rightarrow \infty} 0 .$$

## Other Results

- ▶ Equivalent to do RLHF if  $K \rightarrow \infty$
- ▶ Can be extended with a mix of real and synthetic data

# Self-Consuming Loop

$$p_{t+1} = \arg \max_{p \in \mathcal{P}} \mathbb{E}_{x \sim \hat{p}_{\text{data}}} [\log p(x)] + \lambda \cdot \mathbb{E}_{\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_K \sim p_t} [\log p(\hat{\mathbf{x}})]$$

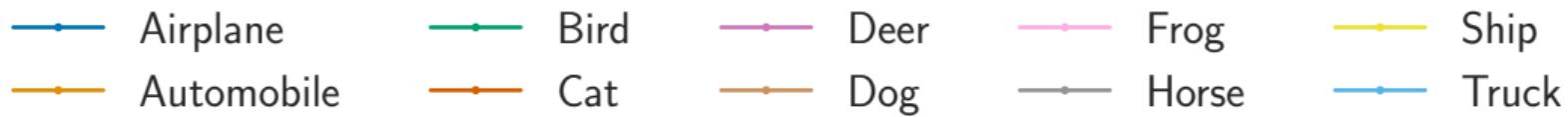
Q: What will happen?

## Results

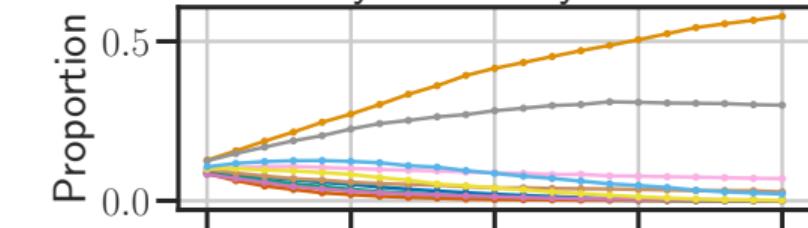
$$\mathbb{E}_{p_t} [e^{r(x)}] \xrightarrow{t \rightarrow \infty} e^{r^*} \quad \text{and} \quad \text{Var}_{p_t} [e^{r(x)}] \xrightarrow{t \rightarrow \infty} 0 .$$

## Other Results

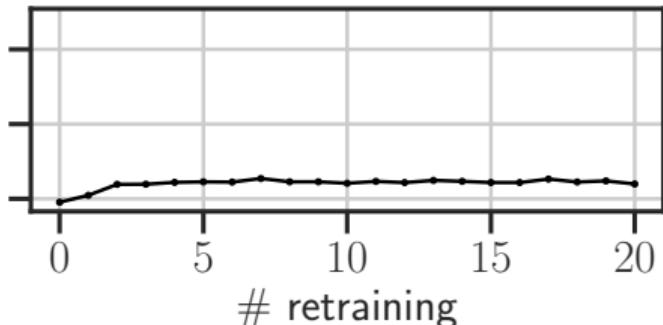
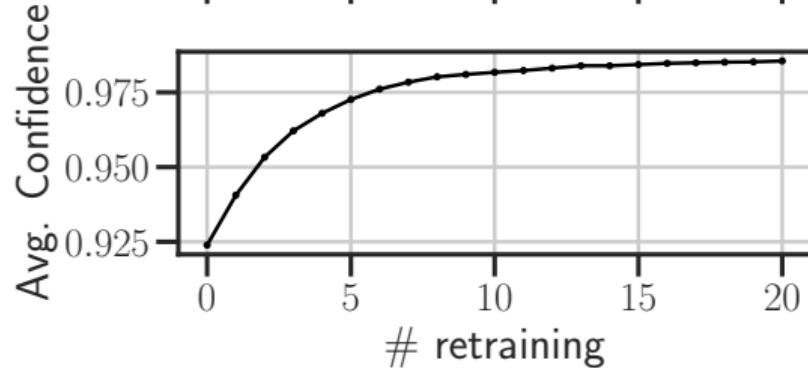
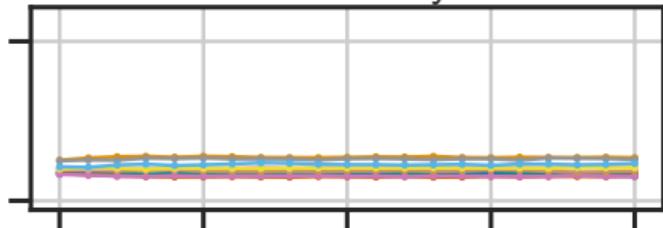
- ▶ Equivalent to do RLHF if  $K \rightarrow \infty$
- ▶ Can be extended with a mix of real and synthetic data



Only Curated Synthetic



Real & Curated Synthetic



- ▶  $\text{reward}(x) = \text{confidence of a pretrained classifier for the image } x$
- ▶  $\lambda = 1/2$

# Conclusion and Future Work

## Future Work

- ▶ Filtering without Human-Feedback?
  - ↪ Score per sample? Downstream-task specific?

---

<sup>a</sup>M. Jiralerspong et al. "Feature Likelihood Score: Evaluating Generalization of Generative Models Using Samples". In: *NeurIPS* (2023).

<sup>b</sup>R. A. Hemmat et al. "Feedback-guided Data Synthesis for Imbalanced Classification". In: *TMLR* (2024).

<sup>c</sup>T. Karras et al. "Guiding a Diffusion Model with a Bad Version of Itself". In: *arXiv preprint arXiv:2406.02507* (2024).

<sup>d</sup>S. Alemdohammad et al. "Self-Improving Diffusion Models with Synthetic Data". In: *arXiv preprint arXiv:2408.16333* (2024).

# Conclusion and Future Work

## Future Work

- ▶ Filtering without Human-Feedback?
  - ↪ Score per sample? Downstream-task specific?
    - ↪ Feature Likelihood Score (FLS)<sup>a</sup>
    - ↪ Classifier to score the samples

---

<sup>a</sup>M. Jiralerspong et al. "Feature Likelihood Score: Evaluating Generalization of Generative Models Using Samples". In: *NeurIPS* (2023).

<sup>b</sup>R. A. Hemmat et al. "Feedback-guided Data Synthesis for Imbalanced Classification". In: *TMLR* (2024).

<sup>c</sup>T. Karras et al. "Guiding a Diffusion Model with a Bad Version of Itself". In: *arXiv preprint arXiv:2406.02507* (2024).

<sup>d</sup>S. Alemdohammad et al. "Self-Improving Diffusion Models with Synthetic Data". In: *arXiv preprint arXiv:2408.16333* (2024).

# Conclusion and Future Work

## Future Work

- ▶ Filtering without Human-Feedback?
  - ↪ Score per sample? Downstream-task specific?
    - ↪ Feature Likelihood Score (FLS)<sup>a</sup>
    - ↪ Classifier to score the samples
    - ↪ Correlation between accuracy and sample quality?<sup>b</sup>

---

<sup>a</sup>M. Jiralerspong et al. "Feature Likelihood Score: Evaluating Generalization of Generative Models Using Samples". In: *NeurIPS* (2023).

<sup>b</sup>R. A. Hemmat et al. "Feedback-guided Data Synthesis for Imbalanced Classification". In: *TMLR* (2024).

<sup>c</sup>T. Karras et al. "Guiding a Diffusion Model with a Bad Version of Itself". In: *arXiv preprint arXiv:2406.02507* (2024).

<sup>d</sup>S. Alemdohammad et al. "Self-Improving Diffusion Models with Synthetic Data". In: *arXiv preprint arXiv:2408.16333* (2024).

# Conclusion and Future Work

## Future Work

- ▶ Filtering without Human-Feedback?
  - ↪ Score per sample? Downstream-task specific?
    - ↪ Feature Likelihood Score (FLS)<sup>a</sup>
    - ↪ Classifier to score the samples
    - ↪ Correlation between accuracy and sample quality?<sup>b</sup>
  - ↪ Score per distribution?
    - ↪ Computationally intensive

---

<sup>a</sup>M. Jiralerspong et al. "Feature Likelihood Score: Evaluating Generalization of Generative Models Using Samples". In: *NeurIPS* (2023).

<sup>b</sup>R. A. Hemmat et al. "Feedback-guided Data Synthesis for Imbalanced Classification". In: *TMLR* (2024).

<sup>c</sup>T. Karras et al. "Guiding a Diffusion Model with a Bad Version of Itself". In: *arXiv preprint arXiv:2406.02507* (2024).

<sup>d</sup>S. Alejomhammad et al. "Self-Improving Diffusion Models with Synthetic Data". In: *arXiv preprint arXiv:2408.16333* (2024).

# Conclusion and Future Work

## Future Work

- ▶ Filtering without Human-Feedback?
  - ↪ Score per sample? Downstream-task specific?
    - ↪ Feature Likelihood Score (FLS)<sup>a</sup>
    - ↪ Classifier to score the samples
      - ↪ Correlation between accuracy and sample quality?<sup>b</sup>
  - ↪ Score per distribution?
    - ↪ Computationally intensive
  - ↪ Use bad samples/models to improve<sup>cd</sup>

---

<sup>a</sup>M. Jiralerspong et al. "Feature Likelihood Score: Evaluating Generalization of Generative Models Using Samples". In: *NeurIPS* (2023).

<sup>b</sup>R. A. Hemmat et al. "Feedback-guided Data Synthesis for Imbalanced Classification". In: *TMLR* (2024).

<sup>c</sup>T. Karras et al. "Guiding a Diffusion Model with a Bad Version of Itself". In: *arXiv preprint arXiv:2406.02507* (2024).

<sup>d</sup>S. Alemdohammad et al. "Self-Improving Diffusion Models with Synthetic Data". In: *arXiv preprint arXiv:2408.16333* (2024).

# Conclusion and Future Work

## Future Work

- ▶ Filtering without Human-Feedback?
  - ↪ Score per sample? Downstream-task specific?
    - ↪ Feature Likelihood Score (FLS)<sup>a</sup>
    - ↪ Classifier to score the samples
      - ↪ Correlation between accuracy and sample quality?<sup>b</sup>
  - ↪ Score per distribution?
    - ↪ Computationally intensive
  - ↪ Use bad samples/models to improve<sup>cd</sup>

---

<sup>a</sup>M. Jiralerspong et al. "Feature Likelihood Score: Evaluating Generalization of Generative Models Using Samples". In: *NeurIPS* (2023).

<sup>b</sup>R. A. Hemmat et al. "Feedback-guided Data Synthesis for Imbalanced Classification". In: *TMLR* (2024).

<sup>c</sup>T. Karras et al. "Guiding a Diffusion Model with a Bad Version of Itself". In: *arXiv preprint arXiv:2406.02507* (2024).

<sup>d</sup>S. Alejomhammad et al. "Self-Improving Diffusion Models with Synthetic Data". In: *arXiv preprint arXiv:2408.16333* (2024).

# Conclusion and Future Work

## Future Work

- ▶ Filtering without Human-Feedback?
  - ↪ Score per sample? Downstream-task specific?
    - ↪ Feature Likelihood Score (FLS)<sup>a</sup>
    - ↪ Classifier to score the samples
      - ↪ Correlation between accuracy and sample quality?
  - ↪ Score per distribution?
    - ↪ Computationally intensive
  - ↪ Use bad samples/models to improve<sup>cd</sup>

---

<sup>a</sup>M. Jiralerspong et al. "Feature Likelihood Score: Evaluating Generalization of Generative Models Using Samples". In: *NeurIPS* (2023).

<sup>b</sup>R. A. Hemmat et al. "Feedback-guided Data Synthesis for Imbalanced Classification". In: *TMLR* (2024).

<sup>c</sup>T. Karras et al. "Guiding a Diffusion Model with a Bad Version of Itself". In: *arXiv preprint arXiv:2406.02507* (2024).

<sup>d</sup>S. Alejomhammad et al. "Self-Improving Diffusion Models with Synthetic Data". In: *arXiv preprint arXiv:2408.16333* (2024).

Thank You!

- ▶ Alemohammad, S. et al. “Self-Consuming Generative Models Go MAD”. In: *ICLR* (2024).
- ▶ Alemohammad, S. et al. “Self-Improving Diffusion Models with Synthetic Data”. In: *arXiv preprint arXiv:2408.16333* (2024).
- ▶ Bertrand, Q. et al. “On the stability of iterative retraining of generative models on their own data”. In: *ICLR* (2024).
- ▶ Ferbach, D. et al. “Self-Consuming Generative Models with Curated Data Provably Optimize Human Preferences”. In: *NeurIPS* (2024).
- ▶ Gulcehre, C. et al. “Reinforced self-training (REST) for language modeling”. In: (2023). *arXiv: 2308.08998 [cs.CL]*.
- ▶ Hataya, R., H. Bao, and HJ. Arai. “Will Large-scale Generative Models Corrupt Future Datasets?” In: *ICCV*. 2023.
- ▶ Hemmat, R. A. et al. “Feedback-guided Data Synthesis for Imbalanced Classification”. In: *TMLR* (2024).
- ▶ Ho, J., A. Jain, and P. Abbeel. “Denoising diffusion probabilistic models”. In: *NeurIPS* (2020).

- ▶ Jiralerspong, M. et al. “Feature Likelihood Score: Evaluating Generalization of Generative Models Using Samples”. In: *NeurIPS* (2023).
- ▶ Junting, P. et al. “JourneyDB: A Benchmark for Generative Image Understanding”. In: *NeurIPS* (2023).
- ▶ Karras, T. et al. “Guiding a Diffusion Model with a Bad Version of Itself”. In: *arXiv preprint arXiv:2406.02507* (2024).
- ▶ Midjourney. <https://www.midjourney.com/home/>. Version 5.2. Accessed: 2023-09-09. 2023.
- ▶ Pressman, J. D., K. Crowson, and Simulacra Captions Contributors. *Simulacra Aesthetic Captions*. Version 1.0. url <https://github.com/JD-P/simulacra-aesthetic-captions>. Stability AI, 2022.
- ▶ Rafailov, R. et al. “Direct preference optimization: Your language model is secretly a reward model”. In: *NeurIPS* (2024).
- ▶ Schuhmann, C. et al. “Laion-5B: An open large-scale dataset for training next generation image-text models”. In: *NeurIPS* (2022).
- ▶ Shumailov, I. et al. “The Curse of Recursion: Training on Generated Data Makes Models Forget”. In: (2023). arXiv: 2305.17493 [cs.LG].

- ▶ Wang, Z. et al. “Better diffusion models further improve adversarial training”. In: *ICML*. 2023.
- ▶ Ziegler, D. M. et al. “Fine-tuning language models from human preferences”. In: *arXiv preprint arXiv:1909.08593* (2019).