# Classifier-Free Diffusion Guidance

NeurIPS 2021 Workshop

Jonathan Ho, Tim Salimans

Presented by Beomsoon Park

# Contents

C&I LAB

# Introduction

# Introduction

❏ **Fidelity vs Diversity**



Trade-off

⟷

Diversity ↓
Sample Quality ↑

Diversity ↑
Sample Quality ↓

**C&I LAB**

# Introduction

❑ **Fidelity vs Diversity**



**Truncation Trick (GAN)**

# Introduction

❑ **Fidelity vs Diversity**

High Temperature
$\tau > 1$

High Temperature
$\tau < 1$

Low Temperature
Sampling

$\tau$ : Temperature ParameterX

# Introduction

❏ **Truncation-like effect on Diffusion**

| Truncation Trick (GAN) | Low Temperature Sampling | → | **Not applicable directly to the diffusion models** |

**Classifier-guided** Diffusion Model

# Introduction

❏ **Drawbacks of Classifier Guidance**

**Overhead to train the extra classifier**

**Must be trained with noisy data at every training procedure**

**Complicate** Diffusion Model Training Pipeline

**C&I LAB**

# Introduction

❑ **Drawbacks of Classifier Guidance**

---

**Algorithm 1** Classifier guided diffusion sampling, given a diffusion model $(\mu_\theta(x_t), \Sigma_\theta(x_t))$, classifier $p_\phi(y|x_t)$, and gradient scale $s$.

---

Input: class label $y$, gradient scale $s$
$x_T \leftarrow$ sample from $\mathcal{N}(0, \mathbf{I})$
**for all** $t$ from $T$ to 1 **do**
    $\mu, \Sigma \leftarrow \mu_\theta(x_t), \Sigma_\theta(x_t)$
    $x_{t-1} \leftarrow$ sample from $\mathcal{N}(\mu + s\Sigma \nabla_{x_t} \log p_\phi(y|x_t), \Sigma)$
**end for**
**return** $x_0$

---

**Algorithm 2** Classifier guided DDIM sampling, given a diffusion model $\epsilon_\theta(x_t)$, classifier $p_\phi(y|x_t)$, and gradient scale $s$.

---

Input: class label $y$, gradient scale $s$
$x_T \leftarrow$ sample from $\mathcal{N}(0, \mathbf{I})$
**for all** $t$ from $T$ to 1 **do**
    $\hat{\epsilon} \leftarrow \epsilon_\theta(x_t) - \sqrt{1 - \bar{\alpha}_t} \nabla_{x_t} \log p_\phi(y|x_t)$
    $x_{t-1} \leftarrow \sqrt{\bar{\alpha}_{t-1}} \left( \frac{x_t - \sqrt{1-\bar{\alpha}_t}\hat{\epsilon}}{\sqrt{\bar{\alpha}_t}} \right) + \sqrt{1 - \bar{\alpha}_{t-1}}\hat{\epsilon}$
**end for**
**return** $x_0$

Sampling
➡ Mixture of **score estimate** & **classifier gradient**

C&I LAB

9

# Introduction

❑ **Drawbacks of Classifier Guidance**

**Algorithm 1** Classifier guided diffusion sampling, given a diffusion model $(\mu_\theta(x_t), \Sigma_\theta(x_t))$, classifier $p_\phi(y|x_t)$, and gradient scale $s$.

---

Input: class label $y$, gradient scale $s$
$x_T \leftarrow$ sample from $\mathcal{N}(0, \mathbf{I})$
**for all** $t$ from $T$ to 1 **do**
  $\mu, \Sigma \leftarrow \mu_\theta(x_t), \Sigma_\theta(x_t)$
  $x_{t-1} \leftarrow$ sample from $\mathcal{N}(\mu + s\Sigma \nabla_{x_t} \log p_\phi(y|x_t), \Sigma)$
**end for**
**return** $x_0$

---

**Algorithm 2** Classifier guided DDIM sampling, given a diffusion model $\epsilon_\theta(x_t)$, classifier $p_\phi(y|x_t)$, and gradient scale $s$.

---

Input: class label $y$, gradient scale $s$
$x_T \leftarrow$ sample from $\mathcal{N}(0, \mathbf{I})$
**for all** $t$ from $T$ to 1 **do**
  $\hat{\epsilon} \leftarrow \epsilon_\theta(x_t) - \sqrt{1 - \bar{\alpha}_t} \nabla_{x_t} \log p_\phi(y|x_t)$
  $x_{t-1} \leftarrow \sqrt{\bar{\alpha}_{t-1}} \left( \frac{x_t - \sqrt{1 - \bar{\alpha}_t}\hat{\epsilon}}{\sqrt{\bar{\alpha}_t}} \right) + \sqrt{1 - \bar{\alpha}_{t-1}}\hat{\epsilon}$
**end for**
**return** $x_0$

---

Image classifier might be confused as **adversarial attack** on sampling

Whether this method is eligible for classifier-based metrics is **questionable**

C&I LAB

# Introduction

❑ **Classifier-Free Guidance (CFG) Diffusion**

**Not much overhead for training extra classifier**

**Extreme Simplicity** (Only a change of one line)

**Similar FID/IS tradeoff** to that of classifier guidance diffusion

C&I LAB

# **Background**

**C&I LAB**

# Background

❏ **Notation**

Forward Process

$$q(\mathbf{z}_\lambda|\mathbf{x}) = \mathcal{N}(\alpha_\lambda \mathbf{x}, \sigma_\lambda^2 \mathbf{I}), \text{ where } \alpha_\lambda^2 = 1/(1 + e^{-\lambda}), \ \sigma_\lambda^2 = 1 - \alpha_\lambda^2$$

$$q(\mathbf{z}_\lambda|\mathbf{z}_{\lambda'}) = \mathcal{N}((\alpha_\lambda/\alpha_{\lambda'})\mathbf{z}_{\lambda'}, \sigma_{\lambda|\lambda'}^2 \mathbf{I}), \text{ where } \lambda < \lambda', \ \sigma_{\lambda|\lambda'}^2 = (1 - e^{\lambda - \lambda'})\sigma_\lambda^2$$

$$\mathbf{z} = \{\mathbf{z}_\lambda \,|\, \lambda \in [\lambda_{\min}, \lambda_{\max}]\} \qquad \lambda_{\min} < \lambda_{\max} \in \mathbb{R}$$

**Noisy Data**

**Hyperparameter for variance scheduling**
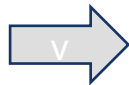
**C&I LAB**

# Background

❑ **Notation**

**Forward Process**

$$q(\mathbf{z}_\lambda|\mathbf{x}) = \mathcal{N}(\alpha_\lambda \mathbf{x}, \sigma_\lambda^2 \mathbf{I}), \text{ where } \alpha_\lambda^2 = 1/(1 + e^{-\lambda}), \ \sigma_\lambda^2 = 1 - \alpha_\lambda^2$$

$$q(\mathbf{z}_\lambda|\mathbf{z}_{\lambda'}) = \mathcal{N}((\alpha_\lambda/\alpha_{\lambda'})\mathbf{z}_{\lambda'}, \sigma_{\lambda|\lambda'}^2 \mathbf{I}), \text{ where } \lambda < \lambda', \ \sigma_{\lambda|\lambda'}^2 = (1 - e^{\lambda - \lambda'})\sigma_\lambda^2$$

$$\lambda = \log(\alpha_\lambda^2/\sigma_\lambda^2)$$
: log SNR of $\boldsymbol{z}_\lambda$

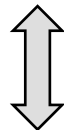⟱ **Formularization of forward process feature**

# Background

❏ **Notation**

**Forward Process**

$$\tilde{\boldsymbol{\mu}}_{\lambda'|\lambda}(\mathbf{z}_\lambda, \mathbf{x}) = e^{\lambda-\lambda'}(\alpha_{\lambda'}/\alpha_\lambda)\mathbf{z}_\lambda + (1-e^{\lambda-\lambda'})\alpha_{\lambda'}\mathbf{x}, \quad \tilde{\sigma}^2_{\lambda'|\lambda} = (1-e^{\lambda-\lambda'})\sigma^2_{\lambda'}$$

$$q(\mathbf{z}_{\lambda'}|\mathbf{z}_\lambda, \mathbf{x}) = \mathcal{N}(\tilde{\boldsymbol{\mu}}_{\lambda'|\lambda}(\mathbf{z}_\lambda, \mathbf{x}), \tilde{\sigma}^2_{\lambda'|\lambda}\mathbf{I})$$

**Final Sample**

$$p_\theta(\mathbf{z}_{\lambda_{\min}}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$

$$p_\theta(\mathbf{z}_{\lambda'}|\mathbf{z}_\lambda) = \mathcal{N}(\tilde{\boldsymbol{\mu}}_{\lambda'|\lambda}(\mathbf{z}_\lambda, \mathbf{x}_\theta(\mathbf{z}_\lambda)), (\tilde{\sigma}^2_{\lambda'|\lambda})^{1-v}(\sigma^2_{\lambda|\lambda'})^v)$$

**Reverse Process**

C&I LAB

# Background

❑ **Notation**

$$p_\theta(\mathbf{z}_{\lambda'}|\mathbf{z}_\lambda) = \mathcal{N}(\tilde{\boldsymbol{\mu}}_{\lambda'|\lambda}(\mathbf{z}_\lambda, \mathbf{x}_\theta(\mathbf{z}_\lambda)), (\tilde{\sigma}^2_{\lambda'|\lambda})^{1-v}(\sigma^2_{\lambda|\lambda'})^v)$$

**DDPM**

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t)) \text{ for } 1 < t \leq T$$

$$\boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t) \begin{cases} \sigma_t^2 = \beta_t \\ \\ \sigma_t^2 = \tilde{\beta}_t = \frac{1-\bar{\alpha}_{t-1}}{1-\bar{\alpha}_t}\beta_t \end{cases}$$

➡ **Similar Results**
(1st option is adopted for this experiment)

$$\underset{\tilde{\beta}_t}{(\tilde{\sigma}^2_{\lambda'|\lambda})^{1-v}} \underset{\beta_t}{(\sigma^2_{\lambda|\lambda'})^v)} : \textbf{log-space linear interpolation}$$

➡ | $v$ : **constant hyperparameter** controlling posterior variance |

# Background

❏ **Training Objective**

$$\mathbb{E}_{\boldsymbol{\epsilon},\lambda}\left[\left\|\boldsymbol{\epsilon}_\theta(\mathbf{z}_\lambda) - \boldsymbol{\epsilon}\right\|_2^2\right]$$ where $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, $\mathbf{z}_\lambda = \alpha_\lambda \mathbf{x} + \sigma_\lambda \boldsymbol{\epsilon}$, $\lambda \sim p(\lambda)$ over $[\lambda_{min}, \lambda_{max}]$

⬇

**Denoising score matching** for all $\lambda$ (noise)

Denoising Score Matching $\qquad \frac{1}{2}\mathbb{E}_{q_\sigma(\tilde{\mathbf{x}}|\mathbf{x})p_{\text{data}}(\mathbf{x})}\left[\left\|\mathbf{s}_\theta(\tilde{\mathbf{x}}) - \nabla_{\tilde{\mathbf{x}}}\log q_\sigma(\tilde{\mathbf{x}}\mid\mathbf{x})\right\|_2^2\right]$

**With this resemblance...** ⇒ $\boldsymbol{\epsilon}_\theta(\mathbf{z}_\lambda) \approx -\sigma_\lambda \nabla_{\mathbf{z}_\lambda}\log p(\mathbf{z}_\lambda)$

**C&I LAB**

17

# Background

❑ **Training Objective**

$$\mathbb{E}_{\boldsymbol{\epsilon},\lambda}\left[\left\|\boldsymbol{\epsilon}_\theta\left(\mathbf{z}_\lambda\right) - \boldsymbol{\epsilon}\right\|_2^2\right] \quad \text{where } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \mathbf{z}_\lambda = \alpha_\lambda\mathbf{x} + \sigma_\lambda\boldsymbol{\epsilon}, \boxed{\lambda \sim p(\lambda) \text{ over } [\lambda_{min}, \lambda_{max}]}$$

$$\boxed{\lambda = -2\log\tan(au + b)}$$

$$a = \arctan\left(e^{-\lambda_{min}/2}\right) - b$$

$$b = \arctan\left(e^{-\lambda_{max}/2}\right)$$

$$u \in [0,1]$$

➡ Sample $\lambda$



Hyperbolic secant distribution

https://en.wikipedia.org/wiki/Hyperbolic_secant_distribution

# Background

❏ **Training Objective**

$$\mathbb{E}_{\boldsymbol{\epsilon},\lambda}\left[\left\|\boldsymbol{\epsilon}_{\theta}(\mathbf{z}_{\lambda}) - \boldsymbol{\epsilon}\right\|_{2}^{2}\right] \quad \text{where } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \mathbf{z}_{\lambda} = \alpha_{\lambda}\mathbf{x} + \sigma_{\lambda}\boldsymbol{\epsilon}, \boxed{\lambda \sim p(\lambda) \text{ over } [\lambda_{min}, \lambda_{max}]}$$



**Weighted** variational lower bound

⬇

**More sophisticated variance scheduling**
➔ Improve sample quality

https://en.wikipedia.org/wiki/Hyperbolic_secant_distribution

# Background

❏ **Sampling Procedure**

$$\epsilon_\theta(\mathbf{z}_\lambda) \approx -\sigma_\lambda \nabla_{\mathbf{z}_\lambda} \log p(\mathbf{z}_\lambda)$$

**Score**

**Learned noise** $\approx$ Estimation of $\nabla_{\mathbf{z}_\lambda} \log p(\mathbf{z}_\lambda)$

⬇

**Correlation** between **sampling from learned diffusion model** and **sampling with Langevin Dynamics**

C&I LAB

# Guidance

C&I LAB

# Guidance

❏ **Classifier Guidance**

**Truncation-like effect** in diffusion models

**Conditional generative modeling**

**Conditions have impacts during training and sampling**

**C&I LAB**

# Guidance

❏ **Classifier Guidance**

$$\epsilon_\theta(\mathbf{z}_\lambda, \mathbf{c}) \approx -\sigma_\lambda \nabla_{\mathbf{z}_\lambda} \log p(\mathbf{z}_\lambda|\mathbf{c})$$

**Diffusion Score during conditional generative modeling**

Parameter determining the strength of classifier

$$\tilde{\epsilon}_\theta(\mathbf{z}_\lambda, \mathbf{c}) = \epsilon_\theta(\mathbf{z}_\lambda, \mathbf{c}) - w\sigma_\lambda \nabla_{\mathbf{z}_\lambda} \log p_\theta(\mathbf{c}|\mathbf{z}_\lambda) \approx -\sigma_\lambda \nabla_{\mathbf{z}_\lambda}[\log p(\mathbf{z}_\lambda|\mathbf{c}) + w \log p_\theta(\mathbf{c}|\mathbf{z}_\lambda)]$$

**classifier**

$$\tilde{p}_\theta(\mathbf{z}_\lambda|\mathbf{c}) \propto p_\theta(\mathbf{z}_\lambda|\mathbf{c})p_\theta(\mathbf{c}|\mathbf{z}_\lambda)^w$$

**Classifier-guided distribution
(model)**

**C&I LAB**

# Guidance

❏ **Classifier Guidance**

$$\tilde{p}_\theta(\mathbf{z}_\lambda | \mathbf{c}) \propto p_\theta(\mathbf{z}_\lambda | \mathbf{c}) \boxed{p_\theta(\mathbf{c} | \mathbf{z}_\lambda)^w}$$

⇨ **Up-weighting the effect of classifier**

⇨ **Higher likelihood to the correct label**
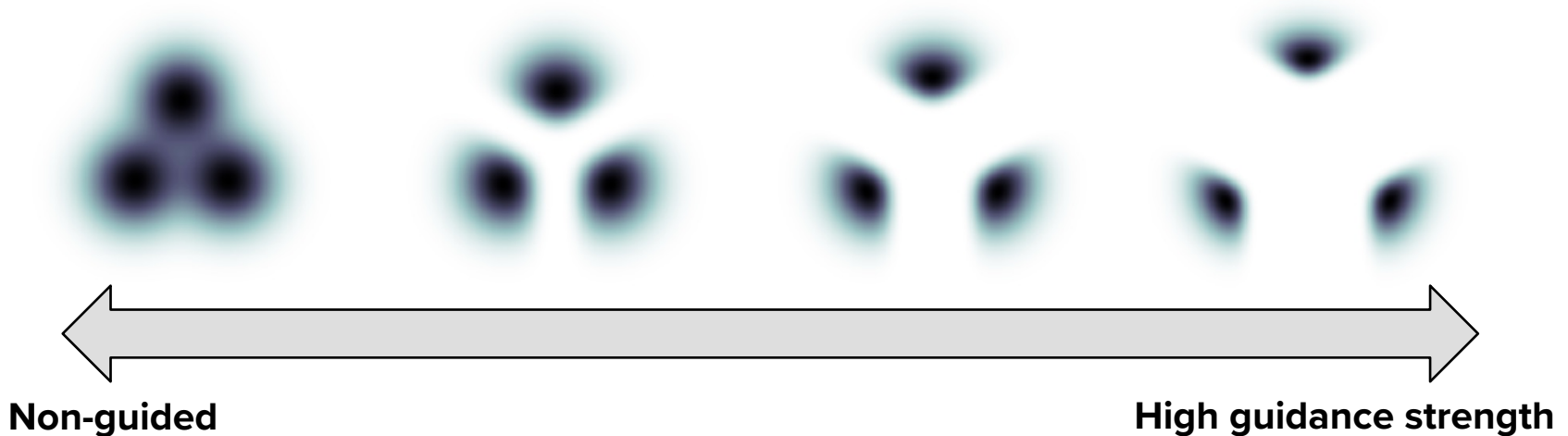
⇨ **Improve IS/FID trade-off by setting $w > 0$**

**C&I LAB**

# Guidance

❑ **Simple Classifier Guidance Experiment**

**Densities of Mixtures of 3 Gaussians**



**Non-guided**

**High guidance strength**

# Guidance

❏ **Classifier Guidance with unconditional model**

**Bayes Rule**

$$p_\theta(\mathbf{z}_\lambda|\mathbf{c})p_\theta(\mathbf{c}|\mathbf{z}_\lambda)^w \propto p_\theta(\mathbf{z}_\lambda)p_\theta(\mathbf{c}|\mathbf{z}_\lambda)^{w+1}$$

$$\boldsymbol{\epsilon}_\theta(\mathbf{z}_\lambda) - (w+1)\sigma_\lambda\nabla_{\mathbf{z}_\lambda}\log p_\theta(\mathbf{c}|\mathbf{z}_\lambda) \approx -\sigma_\lambda\nabla_{\mathbf{z}_\lambda}[\log p(\mathbf{z}_\lambda) + (w+1)\log p_\theta(\mathbf{c}|\mathbf{z}_\lambda)]$$

$$= -\sigma_\lambda\nabla_{\mathbf{z}_\lambda}[\log p(\mathbf{z}_\lambda|\mathbf{c}) + w\log p_\theta(\mathbf{c}|\mathbf{z}_\lambda)]$$

**Same Diffusion score**

**C&I LAB**

# Guidance

❑ **Classifier Guidance**

**Conditional Model**

**Better Performance**

$$\tilde{p}_\theta(\mathbf{z}_\lambda|\mathbf{c}) \propto p_\theta(\mathbf{z}_\lambda|\mathbf{c})p_\theta(\mathbf{c}|\mathbf{z}_\lambda)^w$$

**Unconditional Model**

$$p_\theta(\mathbf{z}_\lambda|\mathbf{c})p_\theta(\mathbf{c}|\mathbf{z}_\lambda)^w \propto p_\theta(\mathbf{z}_\lambda)p_\theta(\mathbf{c}|\mathbf{z}_\lambda)^{w+1}$$

➡ **Classifier-guided conditional model is used for comparison**

# Guidance

❏ **Classifier-Free Guidance (CFG)**

> **Much simpler implementation**

> **Classifier-guided effect without extra classifier**

**C&I LAB**

# Guidance

❏ **Training Classifier-Free Guidance (CFG) diffusion model**

---

**Algorithm 1** Joint training a diffusion model with classifier-free guidance

---

**Require:** $p_{\text{uncond}}$: probability of unconditional training

1: **repeat**
2:     $(\mathbf{x}, \mathbf{c}) \sim p(\mathbf{x}, \mathbf{c})$
3:     $\mathbf{c} \leftarrow \varnothing$ with probability $p_{\text{uncond}}$
4:     $\lambda \sim p(\lambda)$
5:     $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
6:     $\mathbf{z}_\lambda = \alpha_\lambda \mathbf{x} + \sigma_\lambda \boldsymbol{\epsilon}$
7:     Take gradient step on $\nabla_\theta \left\| \boldsymbol{\epsilon}_\theta(\mathbf{z}_\lambda, \mathbf{c}) - \boldsymbol{\epsilon} \right\|^2$
8: **until** converged

**Similar training process to that of DDPM**

---

C&I LAB

# Guidance

❑ **Training Classifier-Free Guidance (CFG) diffusion model**

**Algorithm 1** Joint training a diffusion model with classifier-free guidance

**Require:** $p_{\text{uncond}}$: probability of unconditional training
1: **repeat**
2:     $(\mathbf{x}, \mathbf{c}) \sim p(\mathbf{x}, \mathbf{c})$
3:     $\mathbf{c} \leftarrow \varnothing$ with probability $p_{\text{uncond}}$
4:     $\lambda \sim p(\lambda)$
5:     $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
6:     $\mathbf{z}_\lambda = \alpha_\lambda \mathbf{x} + \sigma_\lambda \boldsymbol{\epsilon}$
7:     Take gradient step on $\nabla_\theta \left\| \boldsymbol{\epsilon}_\theta(\mathbf{z}_\lambda, \mathbf{c}) - \boldsymbol{\epsilon} \right\|^2$
8: **until** converged

**Unconditional model & Conditional model is being trained with 1 neural network**

⬇

**Hyperparameter $p_{uncond}$ decides which model to train during the current iteration**
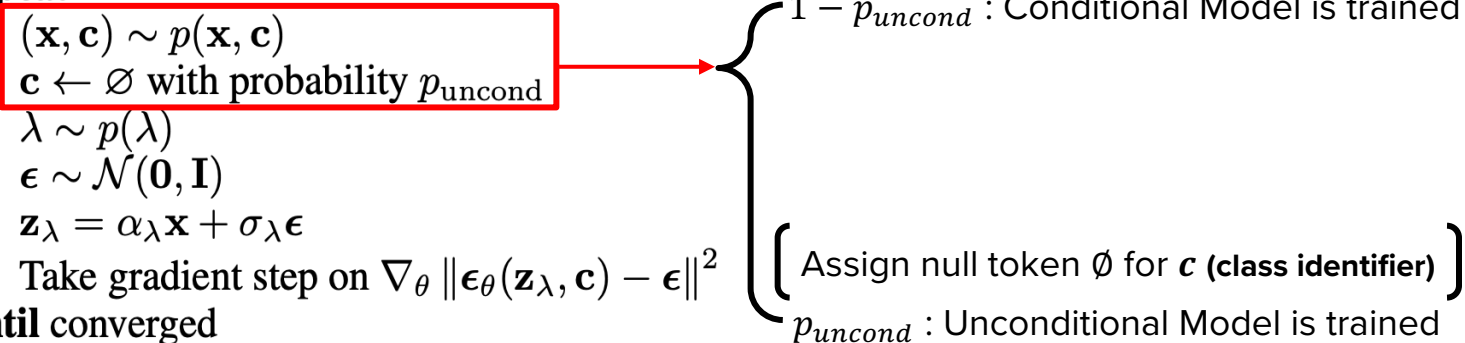
C&I LAB

# Guidance

❑ **Training Classifier-Free Guidance (CFG) diffusion model**

---

**Algorithm 1** Joint training a diffusion model with classifier-free guidance

---

**Require:** $p_{\text{uncond}}$: probability of unconditional training

1: **repeat**
2:     $(\mathbf{x}, \mathbf{c}) \sim p(\mathbf{x}, \mathbf{c})$
3:     $\mathbf{c} \leftarrow \varnothing$ with probability $p_{\text{uncond}}$
4:     $\lambda \sim p(\lambda)$
5:     $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
6:     $\mathbf{z}_\lambda = \alpha_\lambda \mathbf{x} + \sigma_\lambda \boldsymbol{\epsilon}$
7:     Take gradient step on $\nabla_\theta \| \boldsymbol{\epsilon}_\theta(\mathbf{z}_\lambda, \mathbf{c}) - \boldsymbol{\epsilon} \|^2$
8: **until** converged

$1 - p_{uncond}$ : Conditional Model is trained

Assign null token ∅ for $c$ **(class identifier)**

$p_{uncond}$ : Unconditional Model is trained

# Guidance

❑ **Training Classifier-Free Guidance (CFG) diffusion model**

---

**Algorithm 1** Joint training a diffusion model with classifier-free guidance

---

**Require:** $p_{\text{uncond}}$: probability of unconditional training
1: **repeat**
2:     $(\mathbf{x}, \mathbf{c}) \sim p(\mathbf{x}, \mathbf{c})$
3:     $\mathbf{c} \leftarrow \varnothing$ with probability $p_{\text{uncond}}$
4:     $\lambda \sim p(\lambda)$
5:     $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
6:     $\mathbf{z}_\lambda = \alpha_\lambda \mathbf{x} + \sigma_\lambda \boldsymbol{\epsilon}$
7:     Take gradient step on $\nabla_\theta \|\boldsymbol{\epsilon}_\theta(\mathbf{z}_\lambda, \mathbf{c}) - \boldsymbol{\epsilon}\|^2$
8: **until** converged

$$1 - p_{uncond} : \boldsymbol{\epsilon}_\theta(\mathbf{z}_\lambda, \boldsymbol{c})$$

$$p_{uncond} : \boldsymbol{\epsilon}_\theta(\mathbf{z}_\lambda)$$

---

# Guidance

❏ **Sampling from Classifier-Free Guidance (CFG) diffusion model**

**Algorithm 2** Conditional sampling with classifier-free guidance

**Require:** $w$: guidance strength
**Require:** $\mathbf{c}$: conditioning information for conditional sampling ⟹ **Conditions utilized during training**
**Require:** $\lambda_1, \ldots, \lambda_T$: increasing log SNR sequence with $\lambda_1 = \lambda_{\min}, \lambda_T = \lambda_{\max}$

1: $\mathbf{z}_1 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
2: **for** $t = 1, \ldots, T$ **do**
   ▷ Form the classifier-free guided score at log SNR $\lambda_t$
3: $\tilde{\boldsymbol{\epsilon}}_t = (1 + w)\boldsymbol{\epsilon}_\theta(\mathbf{z}_t, \mathbf{c}) - w\boldsymbol{\epsilon}_\theta(\mathbf{z}_t)$ **Main point**
   ▷ Sampling step (could be replaced by another sampler, e.g. DDIM)
4: $\tilde{\mathbf{x}}_t = (\mathbf{z}_t - \sigma_{\lambda_t}\tilde{\boldsymbol{\epsilon}}_t)/\alpha_{\lambda_t}$
5: $\mathbf{z}_{t+1} \sim \mathcal{N}(\tilde{\boldsymbol{\mu}}_{\lambda_{t+1}|\lambda_t}(\mathbf{z}_t, \tilde{\mathbf{x}}_t), (\tilde{\sigma}^2_{\lambda_{t+1}|\lambda_t})^{1-v}(\sigma^2_{\lambda_t|\lambda_{t+1}})^v)$ if $t < T$ else $\mathbf{z}_{t+1} = \tilde{\mathbf{x}}_t$
6: **end for**
7: **return** $\mathbf{z}_{T+1}$

**Sampling Order**

**Also resemble DDPM sampling procedure**

**C&I LAB**

# Guidance

❑ **Suggested score from Classifier-Free Guidance**

$$\tilde{\boldsymbol{\epsilon}}_\theta(\mathbf{z}_\lambda, \mathbf{c}) = (1 + w)\boldsymbol{\epsilon}_\theta(\mathbf{z}_\lambda, \mathbf{c}) - w\boldsymbol{\epsilon}_\theta(\mathbf{z}_\lambda)$$

⬇

$$\tilde{\boldsymbol{\epsilon}}_\theta(\boldsymbol{z}_\lambda, \boldsymbol{c}) = w\left(\boldsymbol{\epsilon}_\theta(\boldsymbol{z}_\lambda, \boldsymbol{c}) - \boldsymbol{\epsilon}_\theta(\boldsymbol{z}_\lambda)\right) + \boldsymbol{\epsilon}_\theta(\boldsymbol{z}_\lambda, \boldsymbol{c})$$

**Parameter** determines
the **strength**

**Implicit classifier guidance**

**C&I LAB**

# Guidance

❏ **Classifier-guided Effect**

$$\left( \boxed{\boldsymbol{\epsilon}_\theta(\mathbf{z}_\lambda, \boldsymbol{c}) - \boldsymbol{\epsilon}_\theta(\mathbf{z}_\lambda)} \right)$$

**Implicit classifier guidance**

⎡ **Recap – Classifier-guided Score** ⎤

$$\tilde{\boldsymbol{\epsilon}}_\theta(\mathbf{z}_\lambda, \mathbf{c}) = \boxed{\boldsymbol{\epsilon}_\theta(\mathbf{z}_\lambda, \mathbf{c})} - w \boxed{\sigma_\lambda \nabla_{\mathbf{z}_\lambda} \log p_\theta(\mathbf{c}|\mathbf{z}_\lambda)}$$

**Score**
**(Conditional generative model)**

**Score**
**(Classifier)**

# Guidance

❑ **Classifier-guided Effect**

$$\left(\boxed{\boldsymbol{\epsilon}_\theta(\boldsymbol{z}_\lambda, \boldsymbol{c}) - \boldsymbol{\epsilon}_\theta(\boldsymbol{z}_\lambda)}\right)$$

**Implicit classifier guidance**

$$\left[ \quad \tilde{\boldsymbol{\epsilon}}_\theta(\boldsymbol{z}_\lambda, \boldsymbol{c}) = w\big(\boldsymbol{\epsilon}_\theta(\boldsymbol{z}_\lambda, \boldsymbol{c}) - \boldsymbol{\epsilon}_\theta(\boldsymbol{z}_\lambda)\big) + \boldsymbol{\epsilon}_\theta(\boldsymbol{z}_\lambda, \boldsymbol{c}) \quad \right]$$

Just estimation, not actual classifier gradient

**Classifier Guidance**
$$\tilde{\boldsymbol{\epsilon}}_\theta(\mathbf{z}_\lambda, \mathbf{c}) = \boxed{\boldsymbol{\epsilon}_\theta(\mathbf{z}_\lambda, \mathbf{c})} - w\boxed{\sigma_\lambda \nabla_{\mathbf{z}_\lambda} \log p_\theta(\mathbf{c}|\mathbf{z}_\lambda)}$$

**Score**
**(Conditional generative model)**

**Score**
**(Classifier)**

**C&I LAB**

# Experiments

C&I LAB

# Experiments

❏ **Experiment**

**Main Purpose**

> Demonstrating attaining IS/FID trade-off similar to that of classifier guidance

⬇

> Same model architecture & hyperparameter settings from classifier guidance

**Suboptimal for classifier-free guidance diffusion model**

**C&I LAB**

# Experiments

❑ **Experiment Settings**

IS/FID score calculation
with 50K Samples

$\lambda_{min} = -20, \lambda_{max} = 20$

**64 X 64**
Conditional
ImageNet
Model

Sampler noise interpolation coefficient $v = \mathbf{0.3}$, **400K** Training steps

**128 X 128**
Conditional
ImageNet
Model

Sampler noise interpolation coefficient $v = \mathbf{0.2}$, **2.7M** Training steps

**C&I LAB**

# Experiments

❑ **Varying Classifier-Free Guidance Strength**

| Model | FID ($\downarrow$) | IS ($\uparrow$) |
|---|---|---|
| ADM (Dhariwal & Nichol, 2021) | 2.07 | - |
| CDM (Ho et al., 2021) | **1.48** | 67.95 |
| Ours | $p_{\mathrm{uncond}} = 0.1/0.2/0.5$ | |
| $w = 0.0$ | 1.8 / 1.8 / 2.21 | 53.71 / 52.9 / 47.61 |
| $w = 0.1$ | 1.55 / 1.62 / 1.91 | 66.11 / 64.58 / 56.1 |
| $w = 0.2$ | 2.04 / 2.1 / 2.08 | 78.91 / 76.99 / 65.6 |
| $w = 0.3$ | 3.03 / 2.93 / 2.65 | 92.8 / 88.64 / 74.92 |
| $w = 0.4$ | 4.3 / 4 / 3.44 | 106.2 / 101.11 / 84.27 |
| $w = 0.5$ | 5.74 / 5.19 / 4.34 | 119.3 / 112.15 / 92.95 |
| $w = 0.6$ | 7.19 / 6.48 / 5.27 | 131.1 / 122.13 / 102 |
| $w = 0.7$ | 8.62 / 7.73 / 6.23 | 141.8 / 131.6 / 109.8 |
| $w = 0.8$ | 10.08 / 8.9 / 7.25 | 151.6 / 140.82 / 116.9 |
| $w = 0.9$ | 11.41 / 10.09 / 8.21 | 161 / 150.26 / 124.6 |
| $w = 1.0$ | 12.6 / 11.21 / 9.13 | 170.1 / 158.29 / 131.1 |
| $w = 2.0$ | 21.03 / 18.79 / 16.16 | 225.5 / 212.98 / 183 |
| $w = 3.0$ | 24.83 / 22.36 / 19.75 | 250.4 / 237.65 / 208.9 |
| $w = 4.0$ | 26.22 / 23.84 / 21.48 | **260.2** / 248.97 / 225.1 |

$w$ ↑

FID ↑ & IS ↑

Table 1: ImageNet 64x64 results ($w = 0.0$ refers to non-guided models).
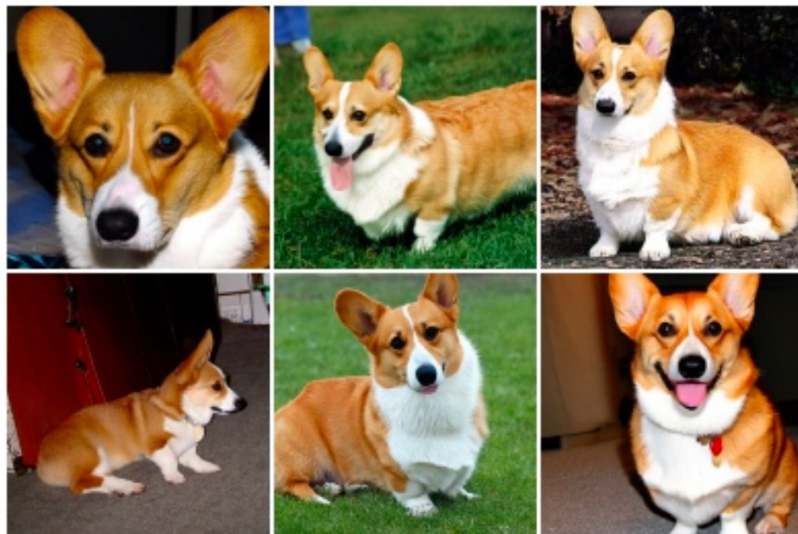
# Experiments

❑ **Varying Classifier-Free Guidance Strength**
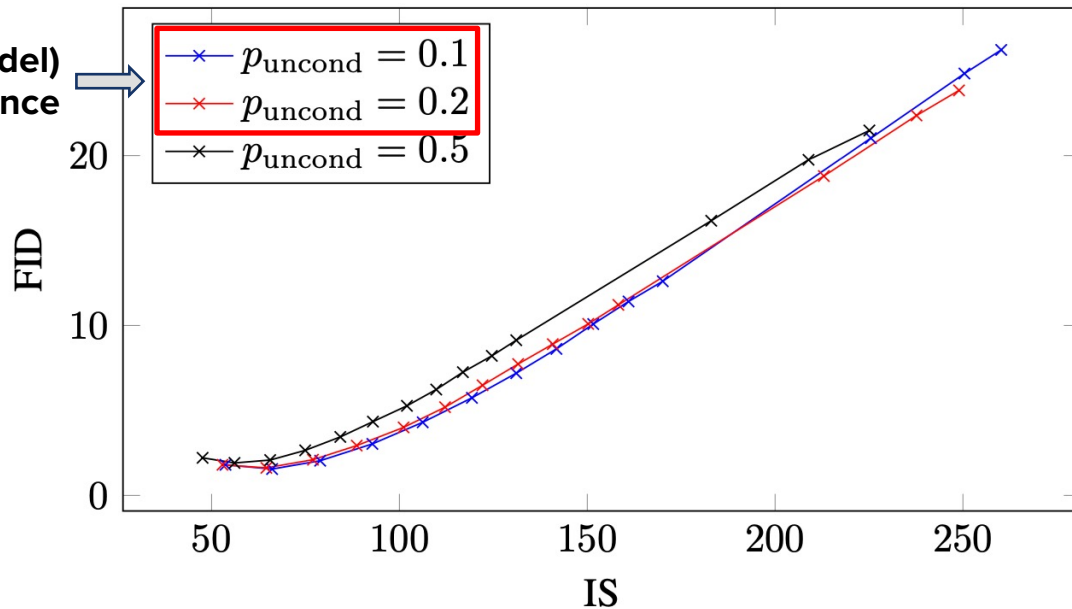
$w = 0$ (Non-guided)  $w = 3.0$

C&I LAB

# Experiments
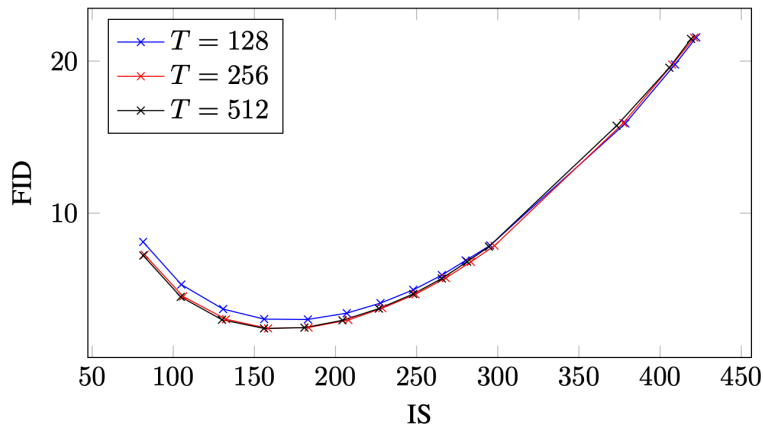
❑ **Varying Unconditional Training Probability**

**(Small portion for unconditional model)**
**Best Performance**

C&I LAB

# Experiments

❑ **Varying Sampling Steps**



| Model | FID (↓) | IS (↑) |
|---|---|---|
| BigGAN-deep, max IS (Brock et al., 2019) | 25 | 253 |
| BigGAN-deep (Brock et al., 2019) | 5.7 | 124.5 |
| CDM (Ho et al., 2021) | 3.52 | 128.8 |
| LOGAN (Wu et al., 2019) | 3.36 | 148.2 |
| ADM-G (Dhariwal & Nichol, 2021) | 2.97 | - |
| Ours | $T = 128/256/1024$ | |
| $w = 0.0$ | 8.11 / 7.27 / 7.22 | 81.46 / 82.45 / 81.54 |
| $w = 0.1$ | 5.31 / 4.53 / 4.5 | 105.01 / 106.12 / 104.67 |
| $w = 0.2$ | 3.7 / 3.03 / 3 | 130.79 / 132.54 / 130.09 |
| $w = 0.3$ | 3.04 / **2.43** / **2.43** | 156.09 / 158.47 / 156 |
| $w = 0.4$ | 3.02 / 2.49 / 2.48 | 183.01 / 183.41 / 180.88 |
| $w = 0.5$ | 3.43 / 2.98 / 2.96 | 206.94 / 207.98 / 204.31 |
| $w = 0.6$ | 4.09 / 3.76 / 3.73 | 227.72 / 228.83 / 226.76 |
| $w = 0.7$ | 4.96 / 4.67 / 4.69 | 247.92 / 249.25 / 247.89 |
| $w = 0.8$ | 5.93 / 5.74 / 5.71 | 265.54 / 267.99 / 265.52 |
| $w = 0.9$ | 6.89 / 6.8 / 6.81 | 280.19 / 283.41 / 281.14 |
| $w = 1.0$ | 7.88 / 7.86 / 7.8 | 295.29 / 297.98 / 294.56 |
| $w = 2.0$ | 15.9 / 15.93 / 15.75 | 378.56 / 377.37 / 373.18 |
| $w = 3.0$ | 19.77 / 19.77 / 19.56 | 409.16 / 407.44 / 405.68 |
| $w = 4.0$ | 21.55 / 21.53 / 21.45 | **422.29** / 421.03 / 419.06 |

Table 2: ImageNet 128x128 results ($w = 0.0$ refers to non-guided models).

C&I LAB

# Experiments

❑ **Varying Sampling Steps**

**Classifier-Guidance Diffusion**

$T = 256$

Best balance
between sample quality & sampling speed

CFG Diffusion model should go through **2 times of forward process**
**(Conditional & Unconditional Model Training)**

**Leading to Slow Sampling Speed**

| Model | FID (↓) | IS (↑) |
|---|---|---|
| BigGAN-deep, max IS (Brock et al., 2019) | 25 | 253 |
| BigGAN-deep (Brock et al., 2019) | 5.7 | 124.5 |
| CDM (Ho et al., 2021) | 3.52 | 128.8 |
| LOGAN (Wu et al., 2019) | 3.36 | 148.2 |
| ADM-G (Dhariwal & Nichol, 2021) | 2.97 | - |
| Ours | $T = 128/256/1024$ | |
| $w = 0.0$ | 8.11 / 7.27 / 7.22 | 81.46 / 82.45 / 81.54 |
| $w = 0.1$ | 5.31 / 4.53 / 4.5 | 105.01 / 106.12 / 104.67 |
| $w = 0.2$ | 3.7 / 3.03 / 3 | 130.79 / 132.54 / 130.09 |
| $w = 0.3$ | 3.04 / **2.43** / 2.43 | 156.09 / 158.47 / 156 |
| $w = 0.4$ | 3.02 / 2.49 / 2.48 | 183.01 / 183.41 / 180.88 |
| $w = 0.5$ | 3.43 / 2.98 / 2.96 | 206.94 / 207.98 / 204.31 |
| $w = 0.6$ | 4.09 / 3.76 / 3.73 | 227.72 / 228.83 / 226.76 |
| $w = 0.7$ | 4.96 / 4.67 / 4.69 | 247.92 / 249.25 / 247.89 |
| $w = 0.8$ | 5.93 / 5.74 / 5.71 | 265.54 / 267.99 / 265.52 |
| $w = 0.9$ | 6.89 / 6.8 / 6.81 | 280.19 / 283.41 / 281.14 |
| $w = 1.0$ | 7.88 / 7.86 / 7.8 | 295.29 / 297.98 / 294.56 |
| $w = 2.0$ | 15.9 / 15.93 / 15.75 | 378.56 / 377.37 / 373.18 |
| $w = 3.0$ | 19.77 / 19.77 / 19.56 | 409.16 / 407.44 / 405.68 |
| $w = 4.0$ | 21.55 / 21.53 / 21.45 | **422.29** / 421.03 / 419.06 |

Table 2: ImageNet 128x128 results ($w = 0.0$ refers to non-guided models).

C&I LAB

# Thank you

**C&I LAB**