



Classifier-free Diffusion Guidance with a Robust Energy-Based Classifier



Advisor Prof. Iacopo Masi Filippo Wang

Computer Science Department

OmnAI Lab, Sapienza, University of Rome, Italy

https://omnai.di.uniroma1.it

Co-Advisor

Dr. Maria Rosaria Briglia







Classifier-free Diffusion Guidance with a Robust Energy-Based Classifier

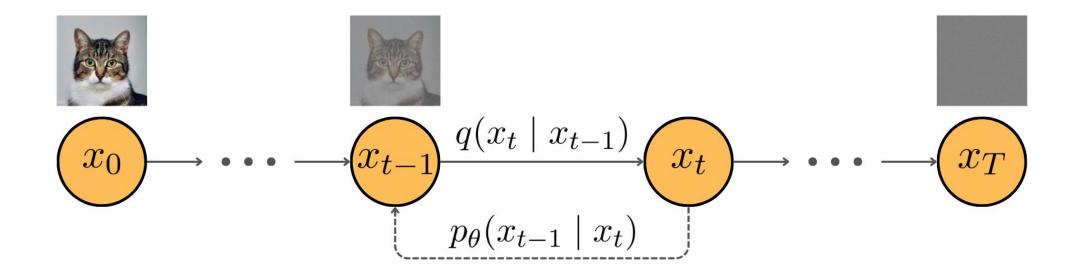
Diffusion Models and Classifier-guidance





Diffusion Models

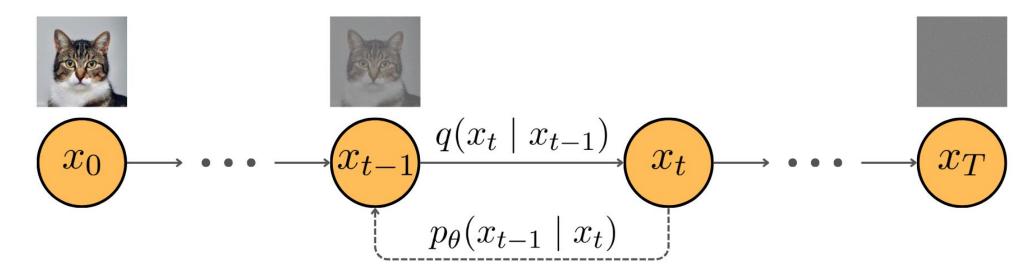


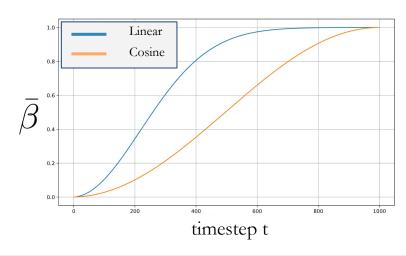




Diffusion Models





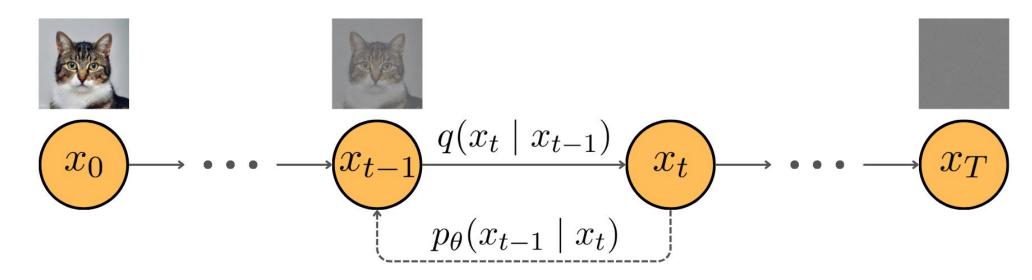


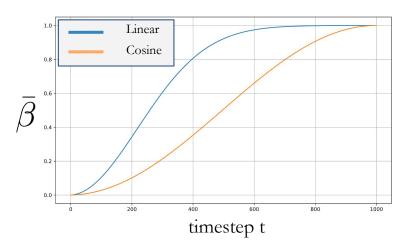
$$q(\mathbf{x}_t \mid \mathbf{x}_0) = \mathcal{N}\left(\mathbf{x}_t; \sqrt{1 - \bar{\beta}_t} \ x_0, \ \bar{\beta}_t \mathbf{I}\right)$$



Diffusion Models







$$q(\mathbf{x}_t \mid \mathbf{x}_0) = \mathcal{N}\left(\mathbf{x}_t; \sqrt{1 - \bar{\beta}_t} \ x_0, \ \bar{\beta}_t \mathbf{I}\right)$$

$$p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}\left(\mathbf{x}_{t-1}; \mu_{\theta}(\mathbf{x}_t, t), \sigma^2\right)$$





1 for all t from T to 1 do

2 |
$$\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$
 ;

$$\begin{array}{c|c}
2 & \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}); \\
3 & \mu \leftarrow \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_{\theta}(\mathbf{x}_t, t) \right) \\
4 & \mathbf{x}_{t-1} \leftarrow \mu + \sigma_t \epsilon;
\end{array}$$

$$\mathbf{x}_{t-1} \leftarrow \mu + \sigma_t \epsilon$$
;

5 return x_0





1 for all t from T to 1 do

2
$$\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$
 :

$$\begin{array}{ll}
2 & \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}); \\
3 & \mu \leftarrow \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_{\theta}(\mathbf{x}_t, t) \right) + s \sigma_t \nabla_{\mathbf{x}_t} \log p_{\phi}(y \mid \mathbf{x}_t) \\
4 & \mathbf{x}_{t-1} \leftarrow \mu + \sigma_t \epsilon;
\end{array}$$

$$\mathbf{x}_{t-1} \leftarrow \mu + \sigma_t \epsilon$$

5 return \mathbf{x}_0





Training is expensive!







Classifier-free Diffusion Guidance with a Robust Energy-Based Classifier

Why Robust Classifiers behave as Generative Models?





Robust Gradients are Perceptually Aligned



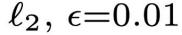
$$\nabla_{\mathbf{x}} \mathcal{L}_{\mathrm{CE}(\mathbf{x},\ y;\ \phi)}$$

Non-Robust

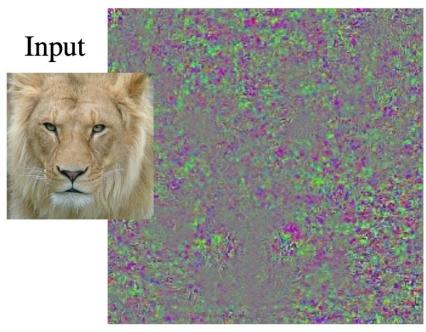
Wang et al. [4]

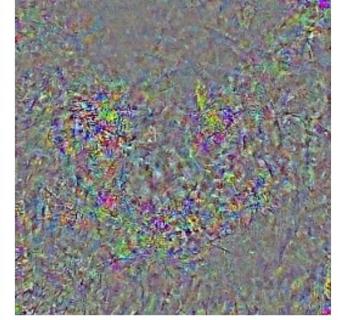
Robust

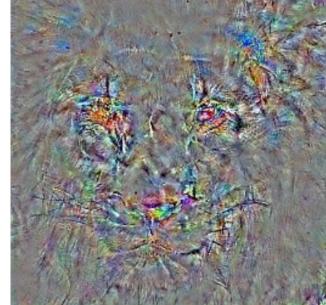
$$\ell_2$$
, ϵ =0.01



$$\ell_2, \; \epsilon = 0.05$$







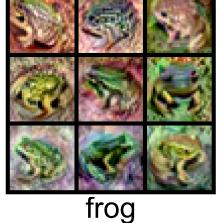


Generating images with a single robust classifier



bird airplane deer cat car











truck





Classifier-free Diffusion Guidance with a Robust Energy-Based Classifier

Robust Classifier Guidance





Naïvely replacing the classifier



airplane

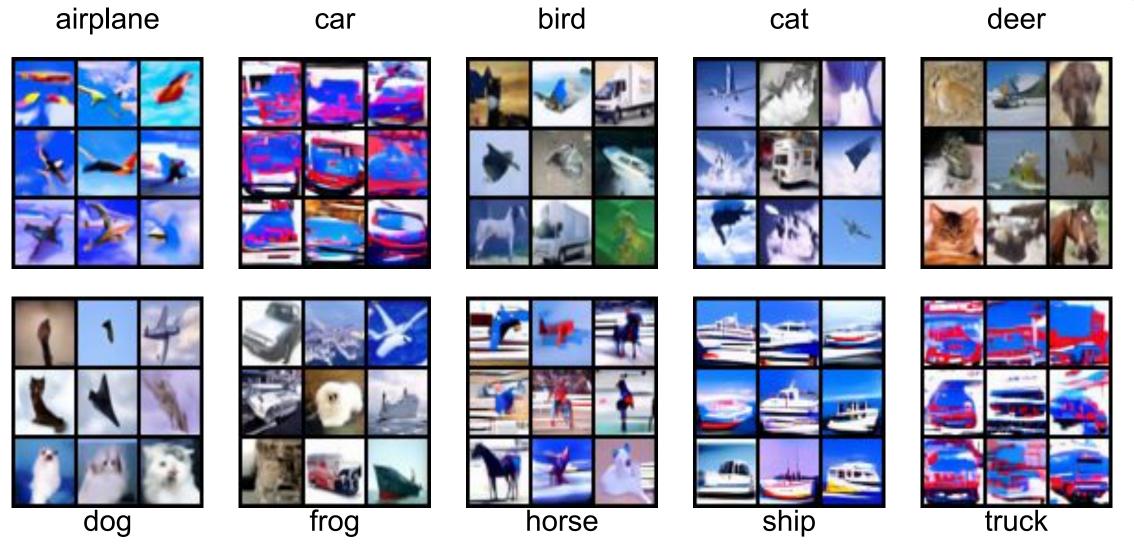






Naïvely replacing the classifier









Classifier-free Diffusion Guidance with a Robust Energy-Based Classifier

Robust Classifier-free guidance

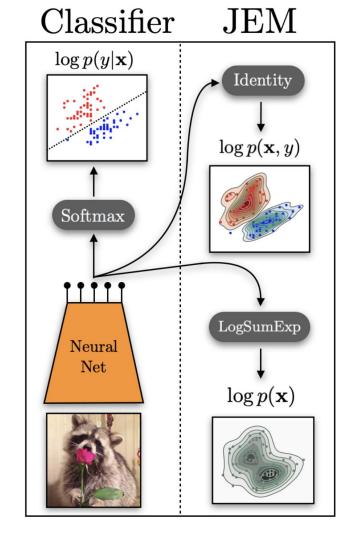




Your classifier is secretly an Energy-Based Model



$$p_{\phi}(y \mid \mathbf{x}) = \frac{p_{\phi}(\mathbf{x}, y)}{p_{\phi}(\mathbf{x})} = \frac{\exp(F_{\phi}(\mathbf{x})[y])}{\sum_{k} \exp(F_{\phi}(\mathbf{x})[y])}$$



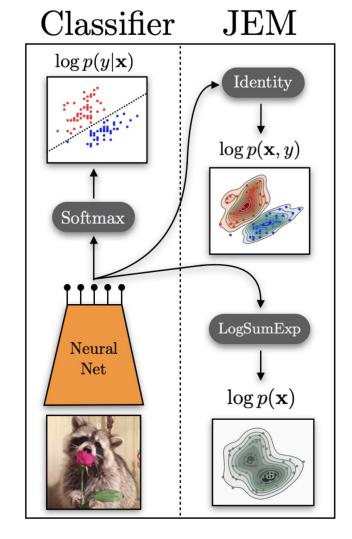


Your classifier is secretly an Energy-Based Model



$$p_{\phi}(y \mid \mathbf{x}) = \frac{p_{\phi}(\mathbf{x}, y)}{p_{\phi}(\mathbf{x})} = \frac{\exp(F_{\phi}(\mathbf{x})[y])}{\sum_{k} \exp(F_{\phi}(\mathbf{x})[y])}$$

$$\log p_{\phi}(y \mid \mathbf{x}) = \log p_{\phi}(\mathbf{x}, y) - \log p_{\phi}(\mathbf{x})$$





Your classifier is secretly an Energy-Based Model



$$p_{\phi}(y \mid \mathbf{x}) = \frac{p_{\phi}(\mathbf{x}, y)}{p_{\phi}(\mathbf{x})} = \frac{\exp(F_{\phi}(\mathbf{x})[y])}{\sum_{k} \exp(F_{\phi}(\mathbf{x})[y])}$$

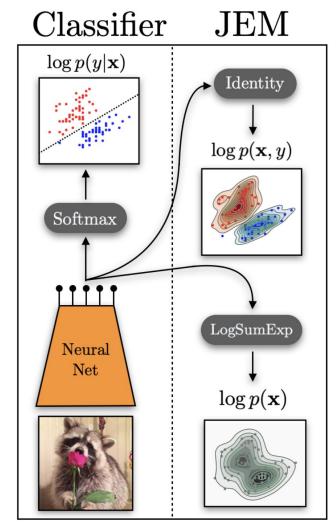
$$\log p_{\phi}(y \mid \mathbf{x}) = \log p_{\phi}(\mathbf{x}, y) - \log p_{\phi}(\mathbf{x})$$

$$= F_{\phi}(\mathbf{x})[y] - \operatorname{LogSumExp}_{y}(F_{\phi}(\mathbf{x})[y])$$

$$= E_{\phi}(\mathbf{x}) - E_{\phi}(\mathbf{x}, y)$$

where
$$E_{\phi}(\mathbf{x}) = -\log p_{\phi}(\mathbf{x})$$

$$E_{\phi}(\mathbf{x}, y) = -\log p_{\phi}(\mathbf{x}, y)$$







1 for all t from T to 1 do

2
$$\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

$$\begin{array}{ll}
2 & \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}); \\
3 & \mu \leftarrow \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_{\theta}(\mathbf{x}_t, t) \right) + s \sigma_t \nabla_{\mathbf{x}_t} \log p_{\phi}(y \mid \mathbf{x}_t) \\
4 & \mathbf{x}_{t-1} \leftarrow \mu + \sigma_t \epsilon;
\end{array}$$

$$\mathbf{x}_{t-1} \leftarrow \mu + \sigma_t \epsilon$$

5 return \mathbf{x}_0





 $-s\hat{\sigma}_t\nabla_{\mathbf{x}_t}E_{\phi}(\mathbf{x},y)$

2 |
$$\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$
 ;

$$\begin{array}{c|c}
c \sim \mathcal{N}(\mathbf{0}, \mathbf{I}); \\
\hline
c \sim \mathcal{N}(\mathbf{0}, \mathbf{I}); \\
c \sim \mathcal{N}(\mathbf{0}, \mathbf{I}); \\
\hline
c \sim \mathcal{N}(\mathbf{0}, \mathbf{I}); \\
c \sim \mathcal{N}(\mathbf{0}, \mathbf{I});$$

$$\mathbf{x}_{t-1} \leftarrow \mu + \sigma_t \epsilon$$

5 return
$$\mathbf{x}_0$$







1 for all t from T to 1 do

$$\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

$$\begin{array}{c|c}
c \sim \mathcal{N}(\mathbf{0}, \mathbf{I}); \\
\mu \leftarrow \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_{\theta}(\mathbf{x}_t, t) \right) - s \hat{\boldsymbol{\sigma}}_t \nabla_{\mathbf{x}_t} E_{\phi}(\mathbf{x}, y) \\
\mathbf{x}_{t-1} \leftarrow \mu + \sigma_t \epsilon;
\end{array}$$

$$\mathbf{x}_{t-1} \leftarrow \mu + \sigma_t \epsilon$$

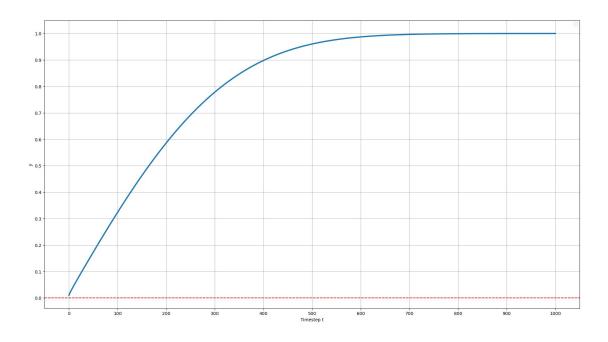
5 return \mathbf{x}_0



Standard guidance scheduler



$$\hat{\sigma}_t = \sqrt{1 - \bar{\alpha}_t}$$

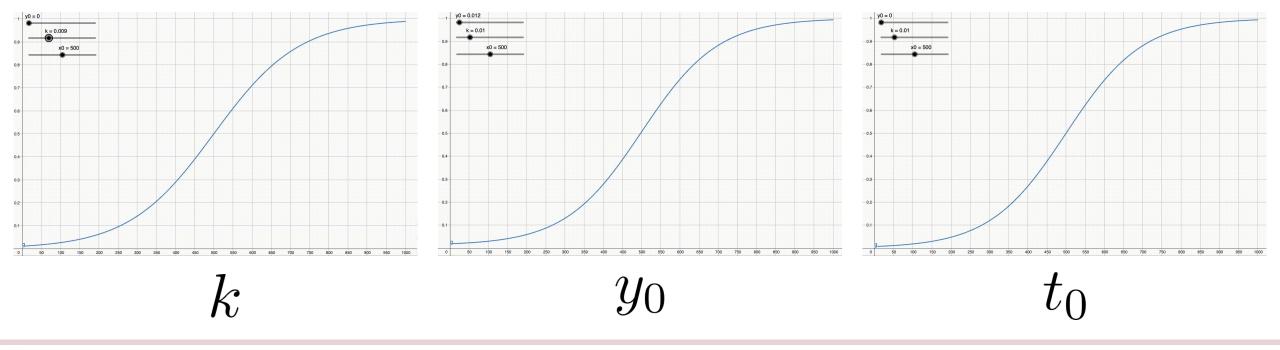




Custom guidance scheduler



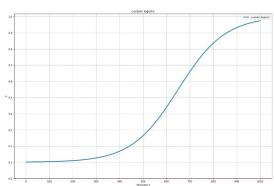
$$\hat{\sigma}_t = (1 - y_0) \cdot \frac{1}{1 + e^{-k(t - t_0)}} + y_0$$





Joint Energy Guidance

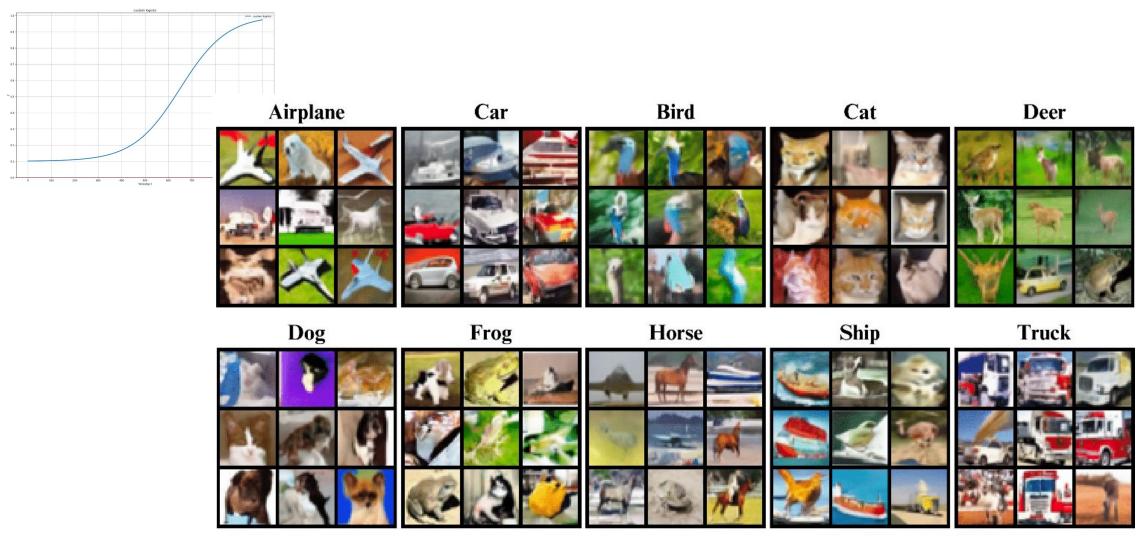






Joint Energy Guidance









Classifier-free Diffusion Guidance with a Robust Energy-Based Classifier

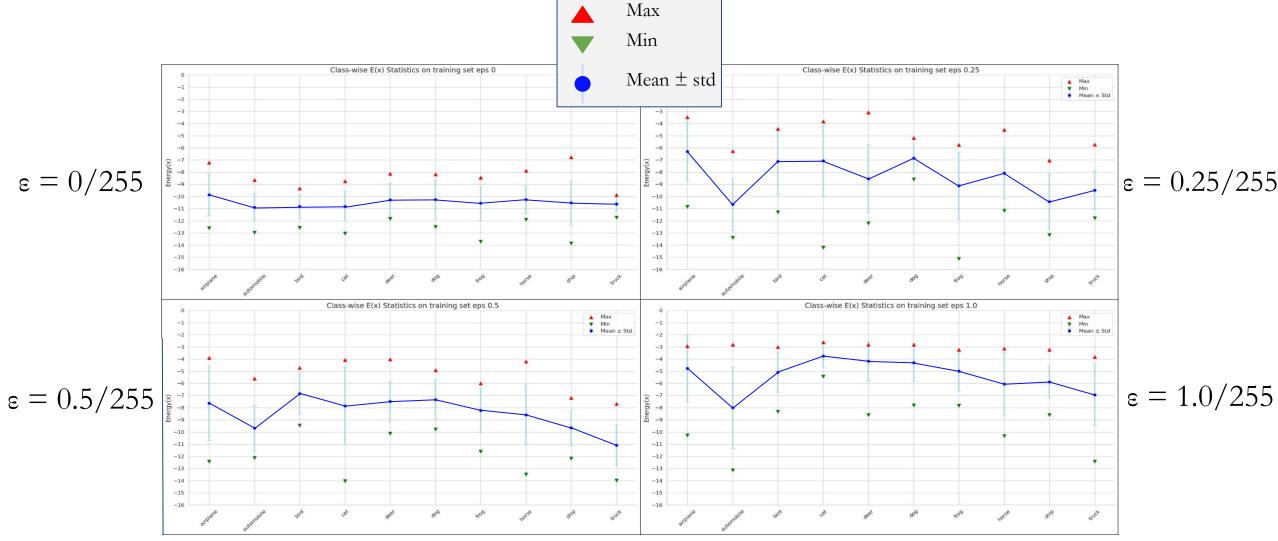
An analysis on training set Energy





E(x) statistics by Class CIFAR-10 training set







Bring back E(x) for guidance



$$\Delta E_{\phi}(\mathbf{x}) = |E_{\phi}(\mathbf{x}^*)[y] - E_{\phi}(\mathbf{x})|$$



Bring back E(x) for guidance



$$\Delta E_{\phi}(\mathbf{x}) = |E_{\phi}(\mathbf{x}^*)[y] - E_{\phi}(\mathbf{x})|$$

$$\mathcal{L} := E_{\phi}(\mathbf{x}, y) + \lambda \Delta E_{\phi}(\mathbf{x})$$





Bring back E(x) for guidance



$$\Delta E_{\phi}(\mathbf{x}) = |E_{\phi}(\mathbf{x}^*)[y] - E_{\phi}(\mathbf{x})|$$

$$\mathcal{L} := E_{\phi}(\mathbf{x}, y) + \lambda \Delta E_{\phi}(\mathbf{x})$$

1 for all t from T to 1 do

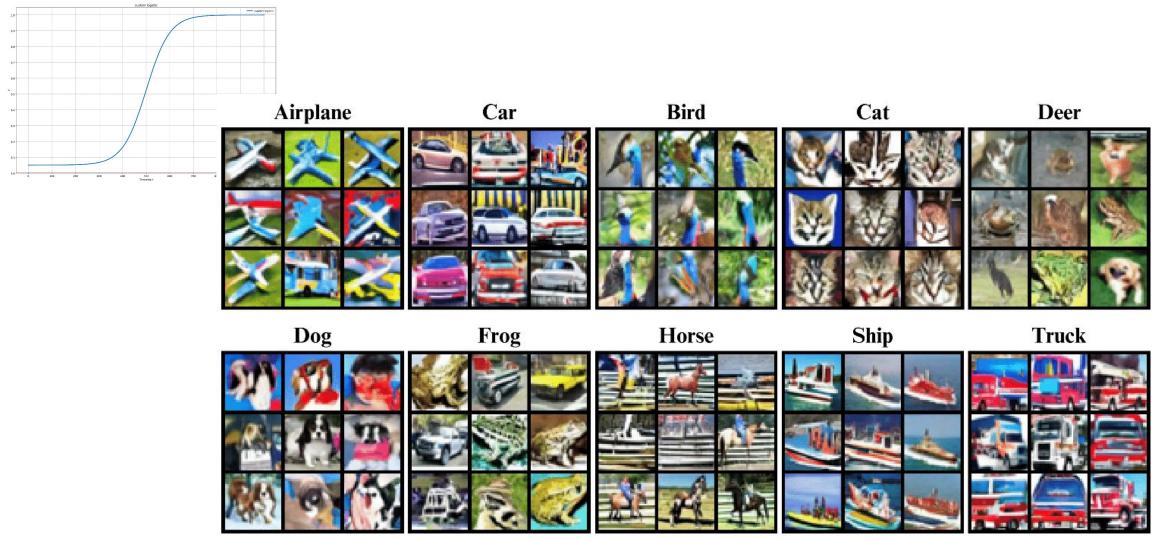
$$\begin{array}{ll}
\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}); & s\hat{\sigma}_t \nabla_{\mathbf{x}_t} \mathcal{L} \\
\mathbf{3} & \mu \leftarrow \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_{\theta}(\mathbf{x}_t, t) \right) - s\hat{\sigma}_t \nabla_{\mathbf{x}_t} E_{\phi}(\mathbf{x}, y) \\
\mathbf{4} & \mathbf{x}_{t-1} \leftarrow \mu + \sigma_t \epsilon;
\end{array}$$

5 return \mathbf{x}_0



Energy is more important than you think







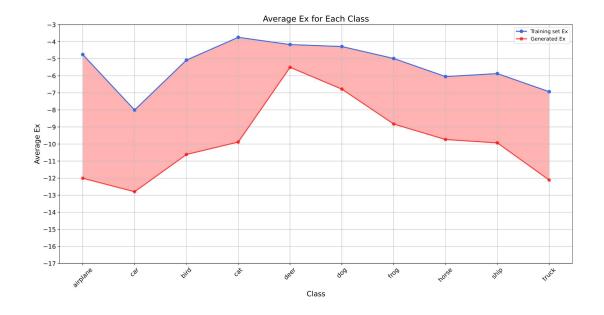
Reducing the Energy Divergence







$abla_{\mathbf{x}_t} \mathcal{L}$







Thanks for listening!

Do you have any questions?



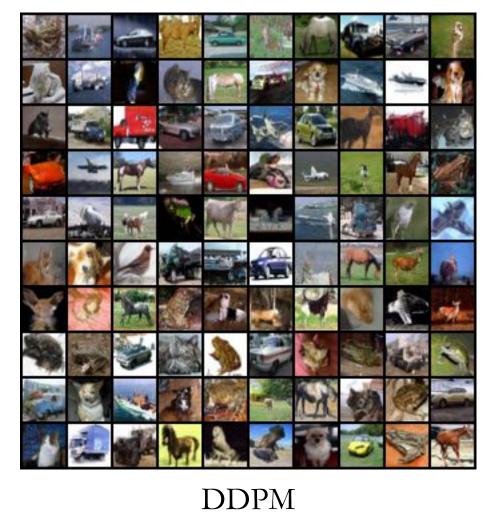






Unconditional Model Quality







DDIM



Diffusion Model Evaluation



	$IS \uparrow$	$FID \downarrow$	Precision ↑	Recall \uparrow
DDPM	8.68	16.49	0.68	0.59
DDIM	8.40	14.79	0.64	0.59

Table 4.2. Evaluation metrics for unconditional guidance using DDPM and DDIM on CIFAR-10 test set, with 1000 and 100 inference steps respectively



Classifier Evaluation



ε -test $\backslash \varepsilon$ -train	0.0	0.25	0.5	1.0
0.0	95.25%	92.77%	90.83%	81.62%
0.25	8.65%	81.21%	$\mathbf{82.40\%}$	75.53%
0.5	0.28%	62.29%	$\boldsymbol{70.17\%}$	68.63%
1.0	0.00%	21.18%	40.48%	$\boldsymbol{52.72\%}$
2.0	0.00%	0.53%	5.22%	18.59%

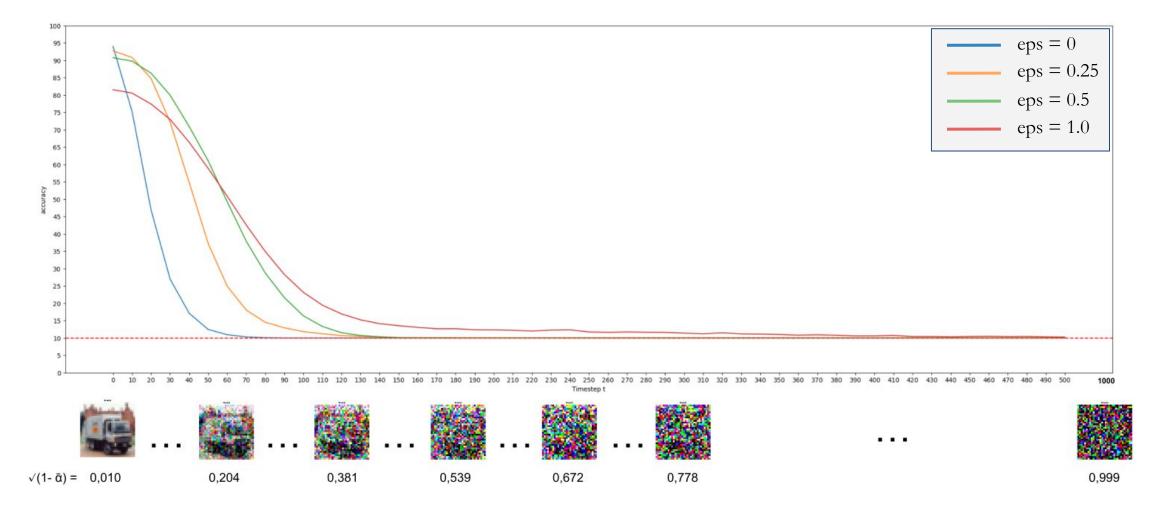
Table 4.1. CIFAR10 L2-norm (ResNet50) 20-steps pgd

test-accuracies on 20-steps PGD-attacks with step size = $2.5 * \frac{\varepsilon_{test}}{20}$



Classifier Accuracy - Gaussian Noise

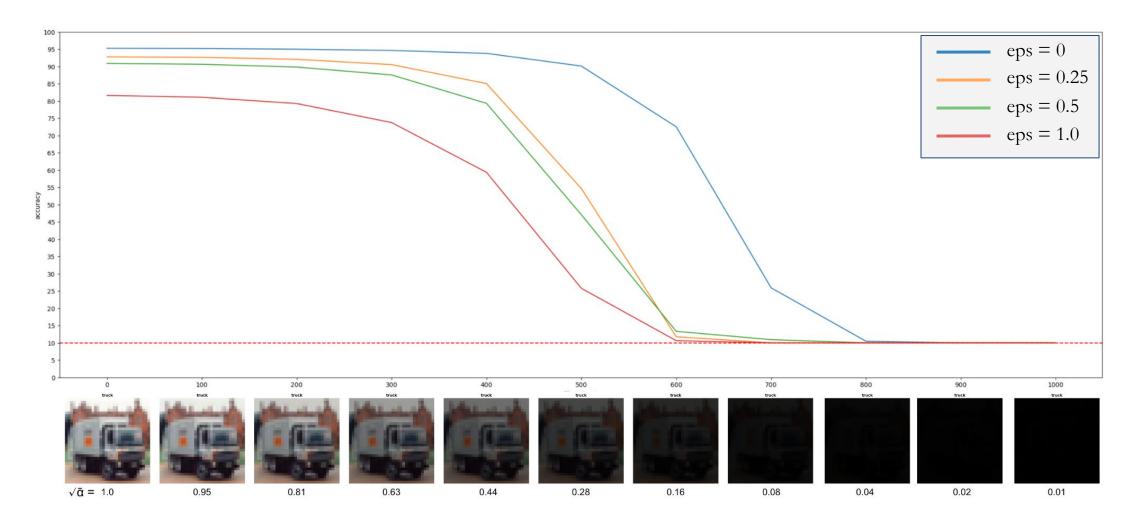






Classifier Accuracy - Darkness









Average E(x, y) and E(x)



