

# Your Diffusion Model is Secretly a Zero-Shot Classifier

Alexander C. Li\*   Mihir Prabhudesai   Shivam Duggal   Ellis Brown  
Deepak Pathak

STRUCT Group Seminar  
Presenter: Zhengbo Xu  
2023.04.02

# OUTLINE

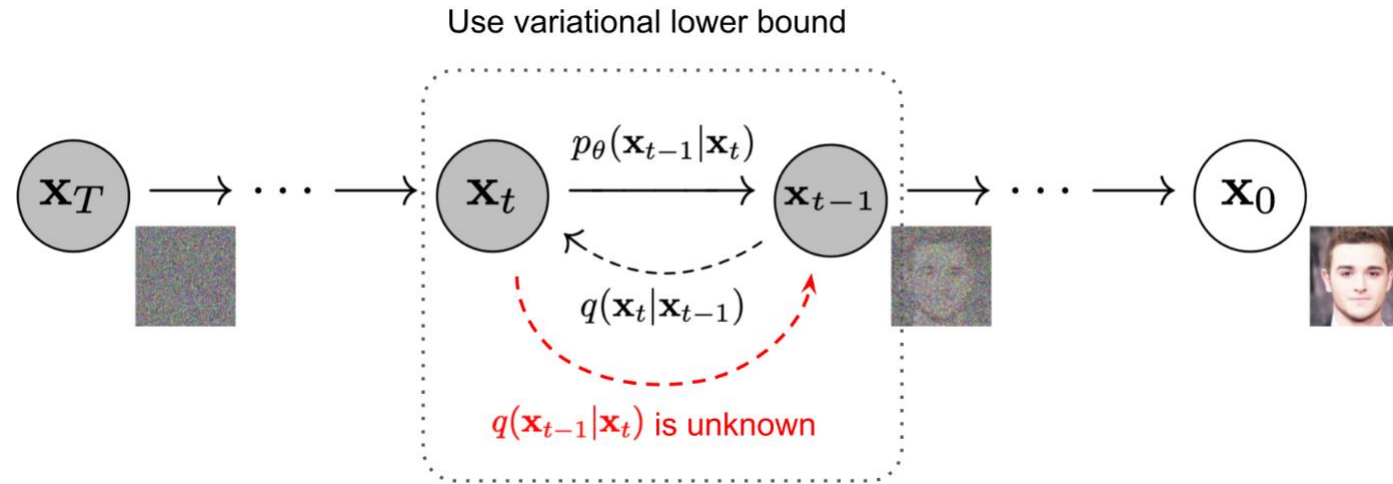
---

- Authorship
- **Background**
- Method
- Experiments
- Conclusion

# BACKGROUND: Diffusion

---

## Overview



$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I}) \quad q(\mathbf{x}_{1:T}|\mathbf{x}_0) = \prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1})$$

# BACKGROUND: Diffusion

---

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}$$

Let  $\alpha_t = 1 - \beta_t$  and  $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}) \quad q(\mathbf{x}_{1:T} | \mathbf{x}_0) = \prod_{t=1}^T q(\mathbf{x}_t | \mathbf{x}_{t-1})$$

The equation can be written as:

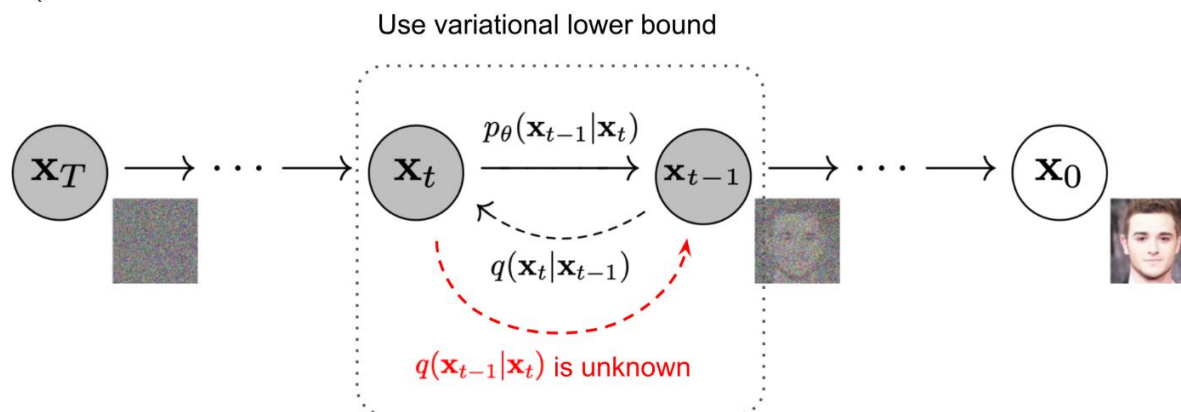
$$\begin{aligned} \mathbf{x}_t &= \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1} \\ &= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_t \alpha_{t-1}} \bar{\boldsymbol{\epsilon}}_{t-2} \\ &= \dots \\ &= \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon} \end{aligned}$$

# BACKGROUND: Diffusion

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}$$

To maximize  $p_\theta(\mathbf{x}_0)$ ,

$$\begin{aligned} -\log p_\theta(\mathbf{x}_0) &\leq -\log p_\theta(\mathbf{x}_0) + D_{\text{KL}}(q(\mathbf{x}_{1:T}|\mathbf{x}_0) \| p_\theta(\mathbf{x}_{1:T}|\mathbf{x}_0)) \\ &= -\log p_\theta(\mathbf{x}_0) + \mathbb{E}_{\mathbf{x}_{1:T} \sim q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[ \log \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{p_\theta(\mathbf{x}_{0:T})/p_\theta(\mathbf{x}_0)} \right] \\ &= -\log p_\theta(\mathbf{x}_0) + \mathbb{E}_q \left[ \log \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{p_\theta(\mathbf{x}_{0:T})} + \log p_\theta(\mathbf{x}_0) \right] \\ &= \mathbb{E}_q \left[ \log \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{p_\theta(\mathbf{x}_{0:T})} \right] \end{aligned}$$



# BACKGROUND: Diffusion

---

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}$$

With a long series of derivation...

$$\mathbb{E}_q \left[ \log \frac{q(\mathbf{x}_{1:T} | \mathbf{x}_0)}{p_\theta(\mathbf{x}_{0:T})} \right] = \mathbb{E}_q \left[ \underbrace{D_{\text{KL}}(q(\mathbf{x}_T | \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_T))}_{L_T} + \sum_{t=2}^T \underbrace{D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t))}_{L_{t-1}} - \underbrace{\log p_\theta(\mathbf{x}_0 | \mathbf{x}_1)}_{L_0} \right]$$

# BACKGROUND: Diffusion

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}$$

With a long series of derivation...

$$\mathbb{E}_q \left[ \log \frac{q(\mathbf{x}_{1:T} | \mathbf{x}_0)}{p_\theta(\mathbf{x}_{0:T})} \right] = \mathbb{E}_q \left[ \underbrace{D_{\text{KL}}(q(\mathbf{x}_T | \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_T))}_{L_T} + \sum_{t=2}^T \underbrace{D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t))}_{L_{t-1}} - \underbrace{\log p_\theta(\mathbf{x}_0 | \mathbf{x}_1)}_{L_0} \right]$$

With another long series of derivation...

$$\begin{aligned} L_t &= \mathbb{E}_{\mathbf{x}_0, \boldsymbol{\epsilon}} \left[ \frac{1}{2 \|\boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t)\|_2^2} \|\tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \mathbf{x}_0) - \boldsymbol{\mu}_\theta(\mathbf{x}_t, t)\|^2 \right] \\ &= \mathbb{E}_{\mathbf{x}_0, \boldsymbol{\epsilon}} \left[ \frac{1}{2 \|\boldsymbol{\Sigma}_\theta\|_2^2} \left\| \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_t \right) - \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \right) \right\|^2 \right] \\ &= \mathbb{E}_{\mathbf{x}_0, \boldsymbol{\epsilon}} \left[ \frac{(1 - \alpha_t)^2}{2 \alpha_t (1 - \bar{\alpha}_t) \|\boldsymbol{\Sigma}_\theta\|_2^2} \|\boldsymbol{\epsilon}_t - \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)\|^2 \right] \\ &= \mathbb{E}_{\mathbf{x}_0, \boldsymbol{\epsilon}} \left[ \frac{(1 - \alpha_t)^2}{2 \alpha_t (1 - \bar{\alpha}_t) \|\boldsymbol{\Sigma}_\theta\|_2^2} \|\boldsymbol{\epsilon}_t - \boldsymbol{\epsilon}_\theta(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}_t, t)\|^2 \right] \end{aligned}$$

# BACKGROUND: Diffusion

---

Simply, we have

$$\begin{aligned} L_t^{\text{simple}} &= \mathbb{E}_{t \sim [1, T], \mathbf{x}_0, \epsilon_t} \left[ \|\epsilon_t - \epsilon_\theta(\mathbf{x}_t, t)\|^2 \right] \\ &= \mathbb{E}_{t \sim [1, T], \mathbf{x}_0, \epsilon_t} \left[ \|\epsilon_t - \epsilon_\theta(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon_t, t)\|^2 \right] \end{aligned}$$

Main steps of diffusion:

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon$$

$$-\log p_\theta(\mathbf{x}_0) \leq \mathbb{E}_q \left[ \log \frac{q(\mathbf{x}_{1:T} | \mathbf{x}_0)}{p_\theta(\mathbf{x}_{0:T})} \right]$$

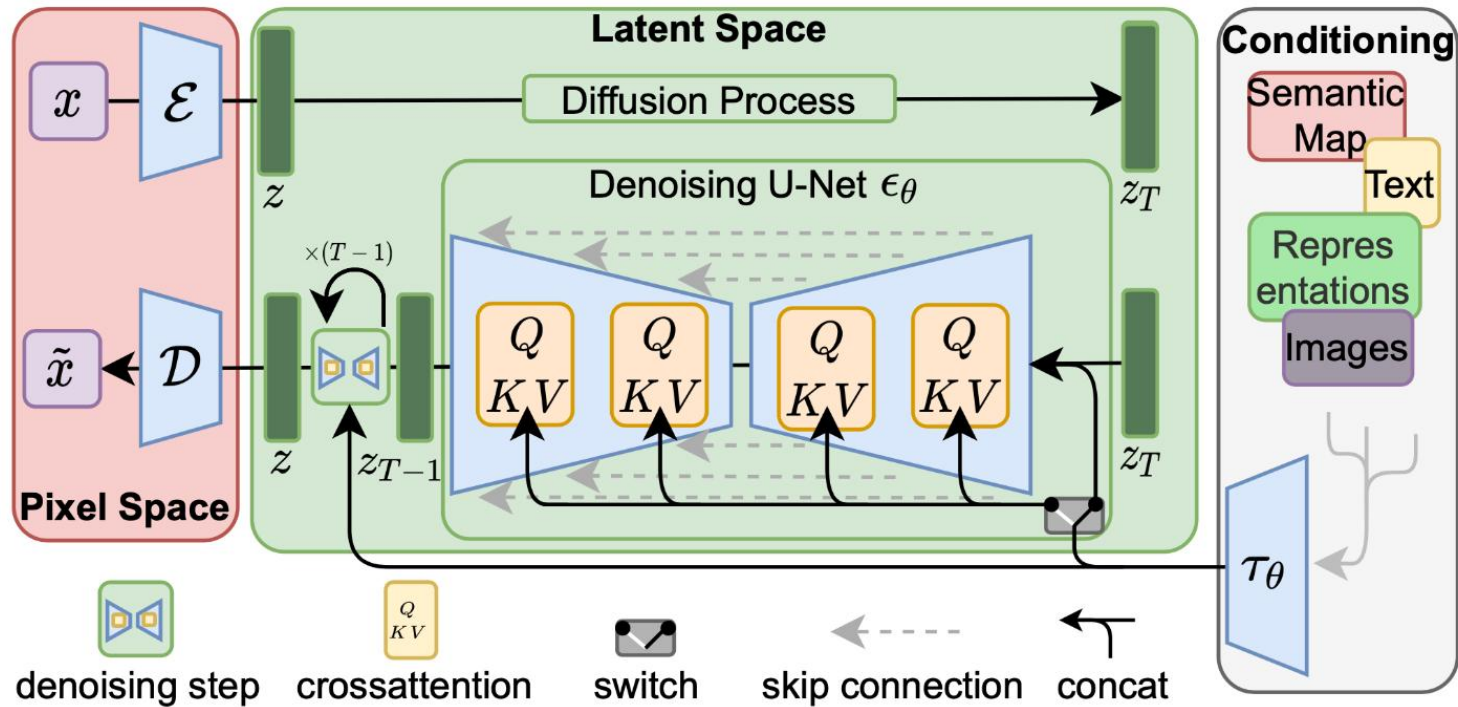
$$\mathbb{E}_q \left[ \log \frac{q(\mathbf{x}_{1:T} | \mathbf{x}_0)}{p_\theta(\mathbf{x}_{0:T})} \right] = \mathbb{E}_q \left[ \underbrace{D_{\text{KL}}(q(\mathbf{x}_T | \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_T))}_{L_T} + \sum_{t=2}^T \underbrace{D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t))}_{L_{t-1}} - \log p_\theta(\mathbf{x}_0 | \mathbf{x}_1) \right]$$

$$L_t^{\text{simple}} = \mathbb{E}_{t \sim [1, T], \mathbf{x}_0, \epsilon_t} \left[ \|\epsilon_t - \epsilon_\theta(\mathbf{x}_t, t)\|^2 \right]$$



# BACKGROUND: Stable Diffusion

## Latent diffusion model



$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d}}\right) \cdot \mathbf{V}$$

$$\text{where } \mathbf{Q} = \mathbf{W}_Q^{(i)} \cdot \varphi_i(\mathbf{z}_i),$$

$$\mathbf{K} = \mathbf{W}_K^{(i)} \cdot \tau_\theta(y),$$

$$\mathbf{V} = \mathbf{W}_V^{(i)} \cdot \tau_\theta(y)$$

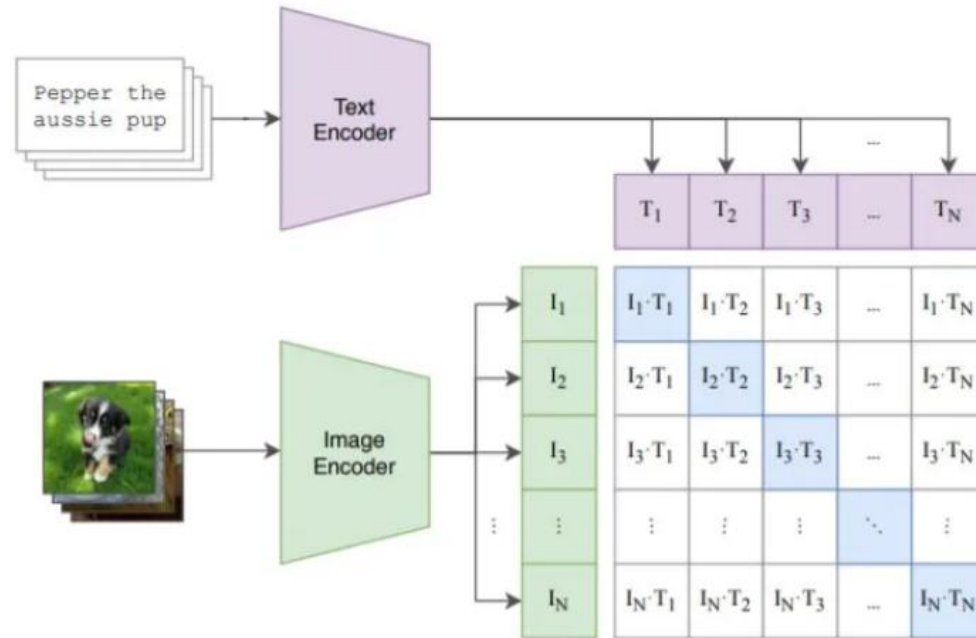
Attention

$$L_{LDM} := \mathbb{E}_{\mathcal{E}(x), \epsilon \sim \mathcal{N}(0,1), t} \left[ \|\epsilon - \epsilon_\theta(z_t, t)\|_2^2 \right]$$

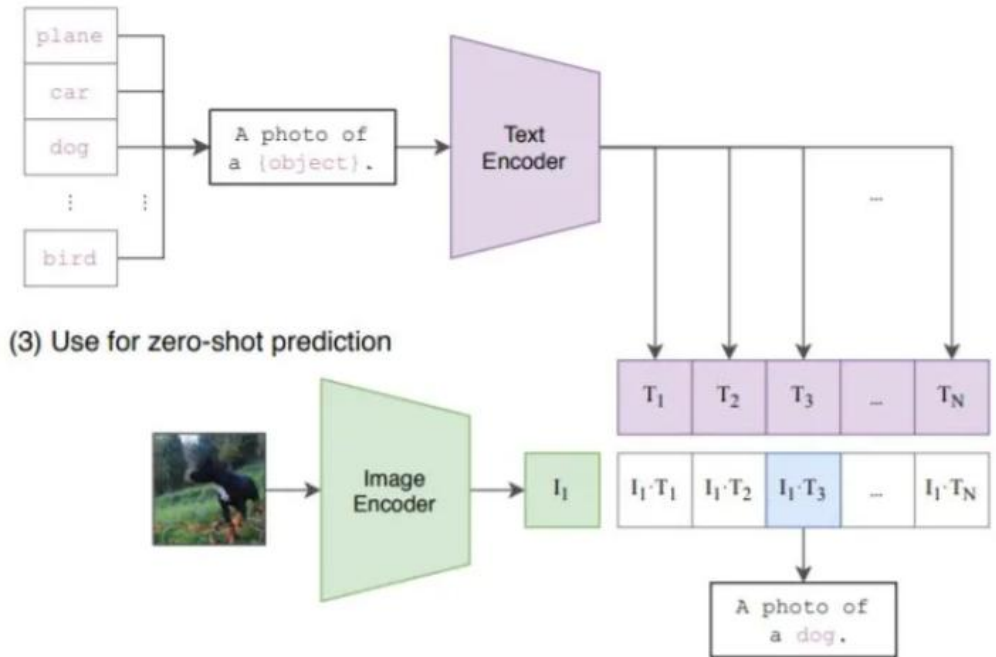
# BACKGROUND: CLIP

## CLIP and OpenCLIP

(1) Contrastive pre-training



(2) Create dataset classifier from label text



# OUTLINE

---

- Authorship
- Background
- Method
- Experiments
- Conclusion

# METHOD

---

How to do classification task using diffusion model?

Generally,

Stable  
Diffusion



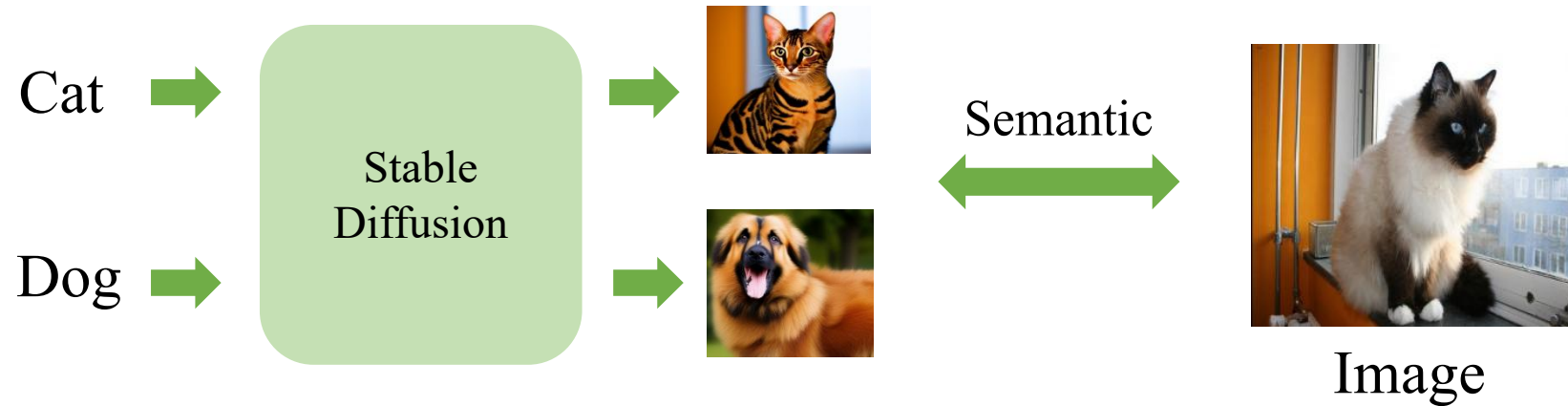
Image

# METHOD

---

How to do classification task using diffusion model?

Generally,



# METHOD

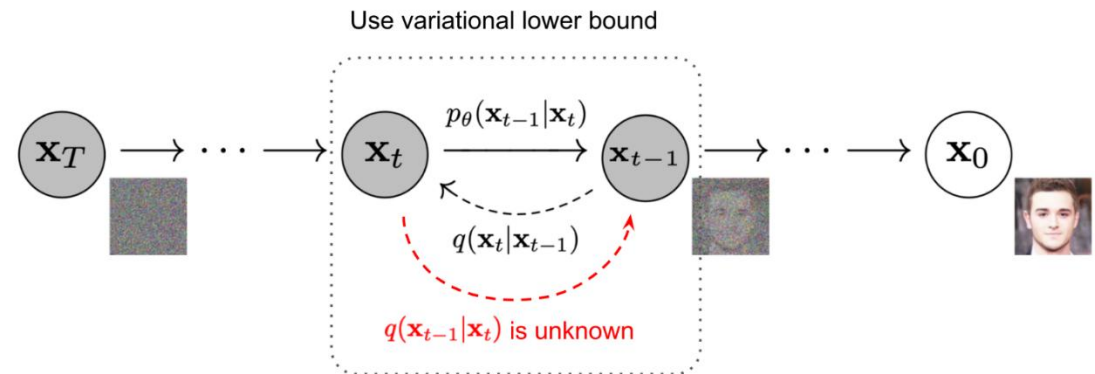
## Quick Review: main steps of diffusion

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}$$

$$-\log p_\theta(\mathbf{x}_0) \leq \mathbb{E}_q \left[ \log \frac{q(\mathbf{x}_{1:T} | \mathbf{x}_0)}{p_\theta(\mathbf{x}_{0:T})} \right]$$

$$\mathbb{E}_q \left[ \log \frac{q(\mathbf{x}_{1:T} | \mathbf{x}_0)}{p_\theta(\mathbf{x}_{0:T})} \right] = \mathbb{E}_q \left[ \underbrace{D_{\text{KL}}(q(\mathbf{x}_T | \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_T))}_{L_T} + \sum_{t=2}^T \underbrace{D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t))}_{L_{t-1}} - \log p_\theta(\mathbf{x}_0 | \mathbf{x}_1) \right]$$

$$L_t^{\text{simple}} = \mathbb{E}_{t \sim [1, T], \mathbf{x}_0, \boldsymbol{\epsilon}_t} \left[ \|\boldsymbol{\epsilon}_t - \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)\|^2 \right]$$



# METHOD

---

How to do classification task using diffusion model?

Theoretically,

$$p_{\theta}(\mathbf{x}_0 | \mathbf{c}) = \int_{\mathbf{x}_{1:T}} p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{c}) d\mathbf{x}_{1:T}$$

# METHOD

---

How to do classification task using diffusion model?

Theoretically,

$$p_{\theta}(\mathbf{x}_0 | \mathbf{c}) = \int_{\mathbf{x}_{1:T}} p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{c}) d\mathbf{x}_{1:T}$$

Bayes' rule:

$$p_{\theta}(\mathbf{c}_i | \mathbf{x}) = \frac{p(\mathbf{c}_i) p_{\theta}(\mathbf{x} | \mathbf{c}_i)}{\sum_j p(\mathbf{c}_j) p_{\theta}(\mathbf{x} | \mathbf{c}_j)}$$



# METHOD

---

How to do classification task using diffusion model?

Theoretically,

$$p_{\theta}(\mathbf{x}_0 | \mathbf{c}) = \int_{\mathbf{x}_{1:T}} p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{c}) d\mathbf{x}_{1:T}$$

Bayes' rule:

$$p_{\theta}(\mathbf{c}_i | \mathbf{x}) = \frac{p(\mathbf{c}_i) p_{\theta}(\mathbf{x} | \mathbf{c}_i)}{\sum_j p(\mathbf{c}_j) p_{\theta}(\mathbf{x} | \mathbf{c}_j)}$$

Suppose  $p(\mathbf{c}_i) = \frac{1}{N}$  ,

We can use  $p_{\theta}(\mathbf{x} | \mathbf{c}_i)$  to predict class.

# METHOD

---

How to predict  $p_\theta(\mathbf{x} | \mathbf{c}_i)$  ?

Using definition is too slow:

$$p_\theta(\mathbf{x}_0 | \mathbf{c}) = \int_{\mathbf{x}_{1:T}} p(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{c}) d\mathbf{x}_{1:T}$$

The variational lower bound:

$$\log p_\theta(\mathbf{x}_0 | \mathbf{c}) \geq \mathbb{E}_q \left[ \log \frac{p_\theta(\mathbf{x}_{0:T}, \mathbf{c})}{q(\mathbf{x}_{1:T} | \mathbf{x}_0)} \right]$$

$$-\log p_\theta(\mathbf{x}_0) \leq \mathbb{E}_q \left[ \log \frac{q(\mathbf{x}_{1:T} | \mathbf{x}_0)}{p_\theta(\mathbf{x}_{0:T})} \right]$$

# METHOD

---

Using

$$\mathbb{E}_q \left[ \log \frac{q(\mathbf{x}_{1:T} | \mathbf{x}_0)}{p_\theta(\mathbf{x}_{0:T})} \right] = \mathbb{E}_q \left[ \underbrace{D_{\text{KL}}(q(\mathbf{x}_T | \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_T))}_{L_T} + \sum_{t=2}^T \underbrace{D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t))}_{L_{t-1}} - \underbrace{\log p_\theta(\mathbf{x}_0 | \mathbf{x}_1)}_{L_0} \right]$$

$$L_t^{\text{simple}} = \mathbb{E}_{t \sim [1, T], \mathbf{x}_0, \epsilon_t} \left[ \|\epsilon_t - \epsilon_\theta(\mathbf{x}_t, t)\|^2 \right]$$

We have

$$\log p_\theta(\mathbf{x}_0 | \mathbf{c}) \geq \mathbb{E}_q \left[ \log \frac{p_\theta(\mathbf{x}_{0:T}, \mathbf{c})}{q(\mathbf{x}_{1:T} | \mathbf{x}_0)} \right] = -\mathbb{E}_{t, \epsilon} \left[ \|\epsilon - \epsilon_\theta(\mathbf{x}_t, \mathbf{c})\|^2 \right] + C$$

# METHOD

---

Combine with Bayes' rule,

$$\begin{aligned} p_{\theta}(\mathbf{c}_i | \mathbf{x}) &\approx \frac{\exp\{-\mathbb{E}_{t,\epsilon}[\|\epsilon - \epsilon_{\theta}(\mathbf{x}_t, \mathbf{c}_i)\|^2] + C\}}{\sum_j \exp\{-\mathbb{E}_{t,\epsilon}[\|\epsilon - \epsilon_{\theta}(\mathbf{x}_t, \mathbf{c}_j)\|^2] + C\}} \\ &= \frac{\exp\{-\mathbb{E}_{t,\epsilon}[\|\epsilon - \epsilon_{\theta}(\mathbf{x}_t, \mathbf{c}_i)\|^2]\}}{\sum_j \exp\{-\mathbb{E}_{t,\epsilon}[\|\epsilon - \epsilon_{\theta}(\mathbf{x}_t, \mathbf{c}_j)\|^2]\}} \end{aligned}$$

How to predict  $\exp\{-\mathbb{E}_{t,\epsilon}[\|\epsilon - \epsilon_{\theta}(\mathbf{x}_t, \mathbf{c}_j)\|^2]\}$  ?

# METHOD

---

Combine with Bayes' rule,

$$\begin{aligned} p_{\theta}(\mathbf{c}_i \mid \mathbf{x}) &\approx \frac{\exp\{-\mathbb{E}_{t,\epsilon}[\|\epsilon - \epsilon_{\theta}(\mathbf{x}_t, \mathbf{c}_i)\|^2] + C\}}{\sum_j \exp\{-\mathbb{E}_{t,\epsilon}[\|\epsilon - \epsilon_{\theta}(\mathbf{x}_t, \mathbf{c}_j)\|^2] + C\}} \\ &= \frac{\exp\{-\mathbb{E}_{t,\epsilon}[\|\epsilon - \epsilon_{\theta}(\mathbf{x}_t, \mathbf{c}_i)\|^2]\}}{\sum_j \exp\{-\mathbb{E}_{t,\epsilon}[\|\epsilon - \epsilon_{\theta}(\mathbf{x}_t, \mathbf{c}_j)\|^2]\}} \end{aligned}$$

How to predict  $\exp\{-\mathbb{E}_{t,\epsilon}[\|\epsilon - \epsilon_{\theta}(\mathbf{x}_t, \mathbf{c}_j)\|^2]\}$  ?

Monte Carlo!

$$\frac{1}{N} \sum_{i=1}^N \left\| \epsilon_i - \epsilon_{\theta}(\sqrt{\bar{\alpha}_{t_i}} \mathbf{x} + \sqrt{1 - \bar{\alpha}_{t_i}} \epsilon_i, \mathbf{c}_j) \right\|^2$$

# METHOD

---

In short,

We just use Bayes' rule and Monte Carlo to classify.

# METHOD

---

In short,

We just use Bayes' rule and Monte Carlo to classify.

Question:

1. Why Monte Carlo can predict  $\exp\{-\mathbb{E}_{t,\epsilon}[\|\epsilon - \epsilon_\theta(\mathbf{x}_t, \mathbf{c}_j)\|^2]\}$  ?
2. Why we can use lower bound?

# METHOD

---

Why Monte Carlo can predict  $\exp\{-\mathbb{E}_{t,\epsilon}[\|\epsilon - \epsilon_\theta(\mathbf{x}_t, \mathbf{c}_j)\|^2]\}$  ?

Indeed, it's hard to predict real value.

But we only need relative value.

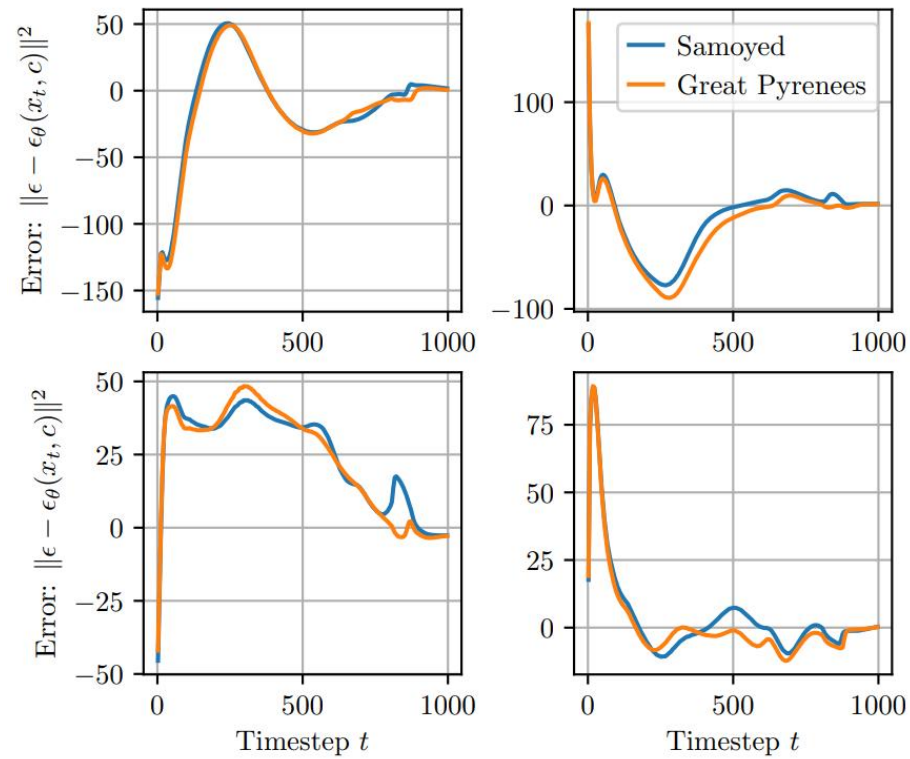
$$\begin{aligned} p_\theta(\mathbf{c}_i | \mathbf{x}) &\approx \frac{\exp\{-\mathbb{E}_{t,\epsilon}[\|\epsilon - \epsilon_\theta(\mathbf{x}_t, \mathbf{c}_i)\|^2] + C\}}{\sum_j \exp\{-\mathbb{E}_{t,\epsilon}[\|\epsilon - \epsilon_\theta(\mathbf{x}_t, \mathbf{c}_j)\|^2] + C\}} \\ &= \frac{\exp\{-\mathbb{E}_{t,\epsilon}[\|\epsilon - \epsilon_\theta(\mathbf{x}_t, \mathbf{c}_i)\|^2]\}}{\sum_j \exp\{-\mathbb{E}_{t,\epsilon}[\|\epsilon - \epsilon_\theta(\mathbf{x}_t, \mathbf{c}_j)\|^2]\}} \\ &= \frac{1}{\sum_j \exp\{\mathbb{E}_{t,\epsilon}[\|\epsilon - \epsilon_\theta(\mathbf{x}_t, \mathbf{c}_i)\|^2 - \|\epsilon - \epsilon_\theta(\mathbf{x}_t, \mathbf{c}_j)\|^2]\}} \end{aligned}$$



# METHOD

---

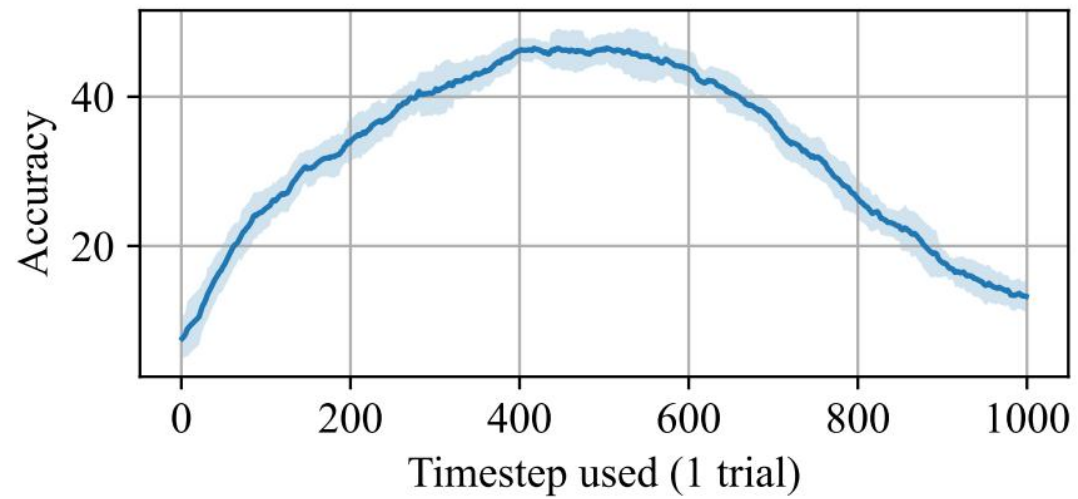
Why we can use lower bound?



# METHOD

---

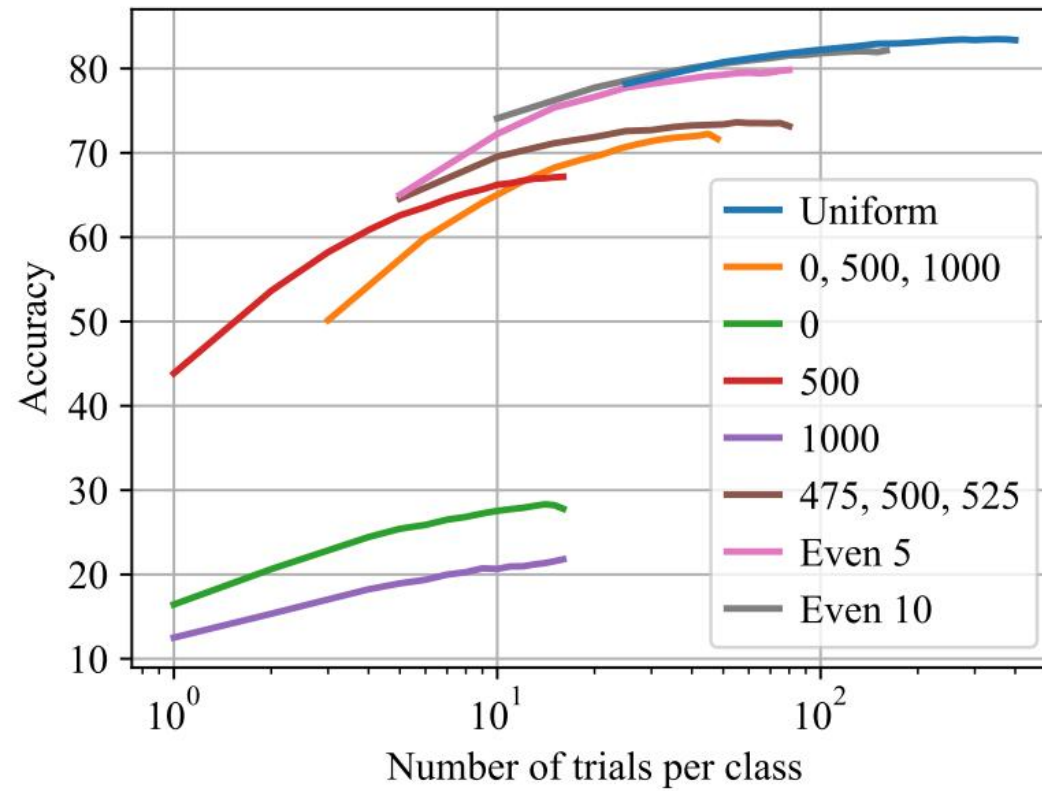
How to choose  $t$  ?



# METHOD

---

How to choose  $t$  ?

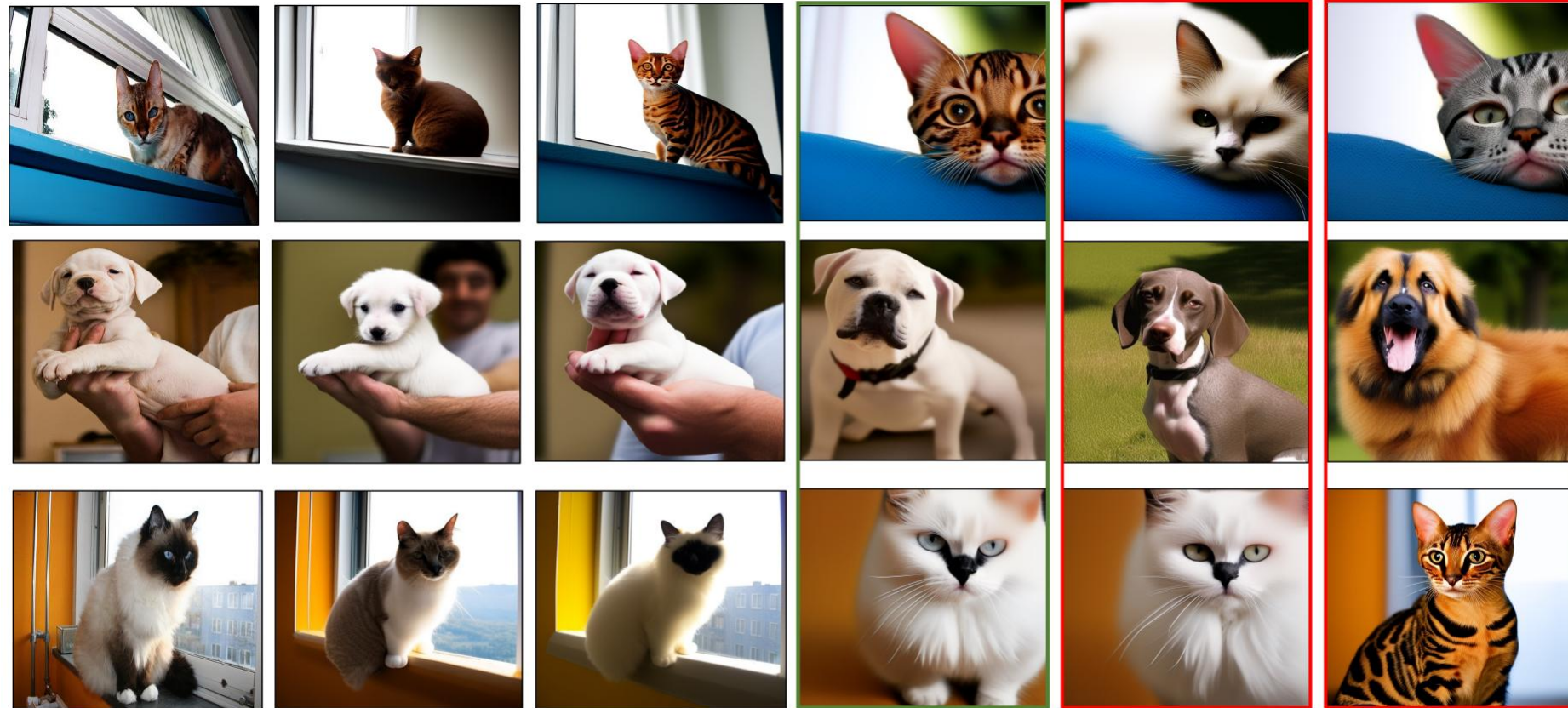


# OUTLINE

---

- Authorship
- Background
- Method
- Experiments
- Conclusion

# EXPERIMENTS: Why does it work?



Input Image

DDIM Inversion  
w/ BLIP caption

DDIM Inversion  
w/ human-modified  
BLIP caption

DDIM Inversion  
w/ **correct** class name  
as prompt

DDIM Inversion  
w/ **incorrect** class  
name as prompt

DDIM Inversion  
w/ **incorrect** class  
name as prompt

# EXPERIMENTS: Zero-Shot Classification

---

	Zero-shot?	Food101	CIFAR10	FGVC	Oxford Pets	Flowers102	MNIST	STL10	ImageNet	ObjectNet
Synthetic SD Data	✓	12.6	35.3	9.4	31.3	22.1	27.9	38.0	18.9	5.2
SD Features	✗	73.0	<b>84.0</b>	<b>35.2</b>	75.9	<b>70.0</b>	<b>98.1</b>	87.2	56.6	10.2
Diffusion Classifier (ours)	✓	<b>77.9</b>	76.3	24.3	<b>85.7</b>	56.8	17.4	<b>94.2</b>	<b>58.4</b>	<b>38.56</b>
CLIP ViT-L/14	✓	93.1	94.5	32.7	93.7	79.3	62.6	99.5	73.5	68.5
OpenCLIP ViT-H/14	✓	92.7	97.3	42.3	94.6	79.9	78.2	98.3	76.8	69.2

Why not as good as CLIP ?

1. A prompt designed for CLIP
2. Training domain without LR, unaesthetic images
3. Choice between log-likelihood and high scores

# EXPERIMENTS: Relational Reasoning

---

## Winoground Benchmark

$\mathbb{I}[\text{score}(C_0, I_0) > \text{score}(C_1, I_0) \text{ AND } \text{score}(C_1, I_1) > \text{score}(C_0, I_1)]$



(a) there is [a mug] in [some grass]



(c) a person [sits] and a dog [stands]



(e) it's a [truck] [fire]



(b) there is [some grass] in [a mug]



(d) a person [stands] and a dog [sits]



(f) it's a [fire] [truck]

*Object*

*Relation*

*Both*



# EXPERIMENTS: Relational Reasoning

## Winoground Benchmark

Model	Object	Relation	Both	Average
Random Chance	25.0	25.0	25.0	25.0
CLIP ViT-L/14	27.0	25.8	57.7	28.2
OpenCLIP ViT-H/14	39.0	<b>26.6</b>	57.7	33.0
Diffusion Classifier (ours)	<b>41.8</b>	25.3	<b>69.2</b>	<b>34.0</b>

✓ Diffusion Classifier ✓ OpenCLIP ✓ CLIP

Mask Mask

"a bird eats a snake" "a snake eats a bird"

✓ Diffusion Classifier ✓ OpenCLIP ✗ CLIP



"there are more ladybugs than flowers" "there are more flowers than ladybugs"

✗ Diffusion Classifier ✗ OpenCLIP ✗ CLIP



"the taller person hugs the shorter person" "the shorter person hugs the taller person"

✓ Diffusion Classifier ✗ OpenCLIP ✗ CLIP



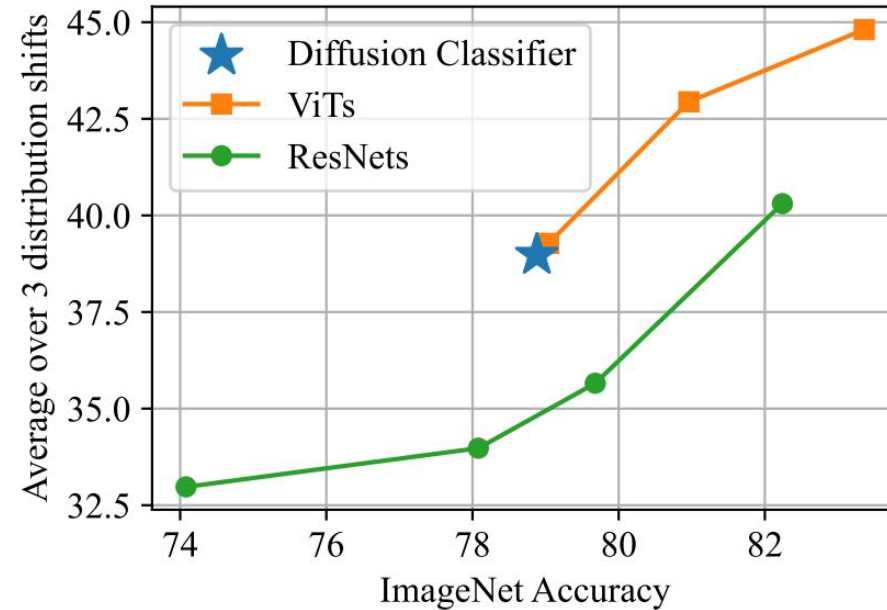
"an old person kisses a young person" "a young person kisses an old person"



# EXPERIMENTS: Supervised Classification

Results with fewer augmentation

Method	ID		OOD	
	IN	IN-v2	IN-A	ObjectNet
ResNet-18	74.1	57.3	15.0	26.6
ResNet-34	78.1	59.8	10.5	31.6
ResNet-50	79.7	61.6	9.8	35.6
ResNet-101	82.2	63.2	19.5	38.2
ViT-L/32	79.0	61.6	26.3	29.9
ViT-L/16	81.0	66.6	25.6	36.7
ViT-B/16	83.4	66.6	30.1	37.8
Diffusion Classifier	78.9	62.1	22.6	32.3



# OUTLINE

---

- Authorship
- Background
- Method
- Experiments
- Conclusion

# CONCLUSION

---

- Diffusion model can be used to zero-shot classification task
- Diffusion model has competitive scores in classification
- Exact choices made during diffusion training affect the classifier

$$\log p_{\theta}(\mathbf{x}_0 \mid \mathbf{c}) \geq \mathbb{E}_q \left[ \log \frac{p_{\theta}(\mathbf{x}_{0:T}, \mathbf{c})}{q(\mathbf{x}_{1:T} \mid \mathbf{x}_0)} \right] = -\mathbb{E}_{t, \epsilon} [\|\epsilon - \epsilon_{\theta}(\mathbf{x}_t, \mathbf{c})\|^2] + C$$

$$\frac{1}{N} \sum_{i=1}^N \left\| \epsilon_i - \epsilon_{\theta}(\sqrt{\bar{\alpha}_t} \mathbf{x} + \sqrt{1 - \bar{\alpha}_t} \epsilon_i, \mathbf{c}_j) \right\|^2$$

Thanks for listening!