

Exploring the Connection between **Robust** and **Generative** Models



github.com/senad96/Robust-Generative

Senad Beadini

Computer Science Department
Sapienza, University of Rome*



Iacopo Masi

Computer Science Department
Sapienza, University of Rome



Exploring the Connection between Robust and Generative Models

Why Robust Models behave as Generatives?



Generative Models



SAPIENZA
UNIVERSITÀ DI ROMA

Generative Models



$$\{\mathbf{x}_i\}$$

Generative Models



$$\{\mathbf{x}_i\} \sim p_{\text{data}}(\mathbf{x})$$

Generative Models



$$\mathbf{x}' \sim p(\mathbf{x}; \boldsymbol{\theta}) \quad \{\mathbf{x}_i\} \sim p_{\text{data}}(\mathbf{X})$$

Generative Models



$$\mathbf{x}' \sim p(\mathbf{x}; \boldsymbol{\theta}) \quad \{\mathbf{x}_i\} \sim p_{\text{data}}(\mathbf{X})$$

GAN

Generative Models



$$\mathbf{x}' \sim p(\mathbf{x}; \boldsymbol{\theta}) \quad \{\mathbf{x}_i\} \sim p_{\text{data}}(\mathbf{X})$$

GAN

VAE

Generative Models



$$\mathbf{x}' \sim p(\mathbf{x}; \boldsymbol{\theta}) \quad \{\mathbf{x}_i\} \sim p_{\text{data}}(\mathbf{X})$$

GAN

VAE

Autoregressive,
Normalizing
Flows,
Invertible NN,
and EBM

Generative Models

$$\mathbf{x}' \sim p(\mathbf{x}; \boldsymbol{\theta}) \quad \{\mathbf{x}_i\} \sim p_{\text{data}}(\mathbf{X})$$

GAN

VAE

Autoregressive,
Normalizing
Flows,
Invertible NN,
and EBM

Diffusion
Models

Generative Models



$$\mathbf{x}' \sim p(\mathbf{x}; \boldsymbol{\theta}) \quad \{\mathbf{x}_i\} \sim p_{\text{data}}(\mathbf{x})$$

GAN

VAE

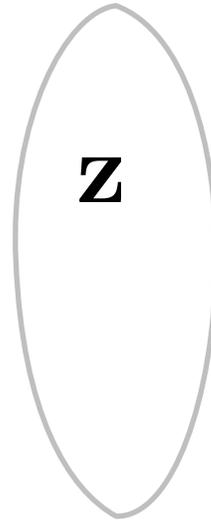
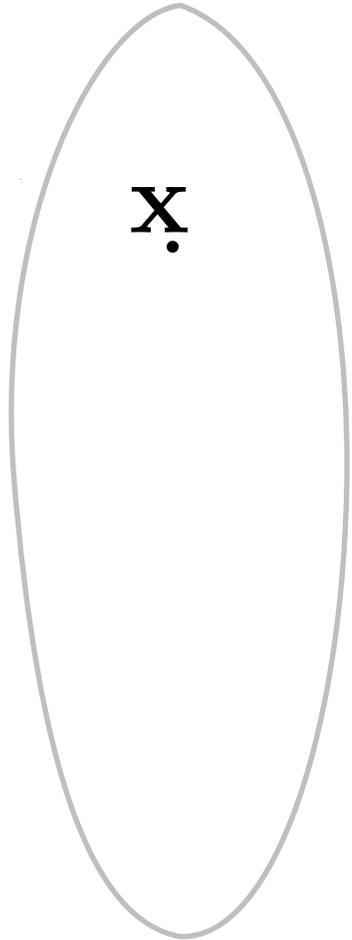
$$p(z|\mathbf{x})$$

Inverting a
discriminative,
robust model

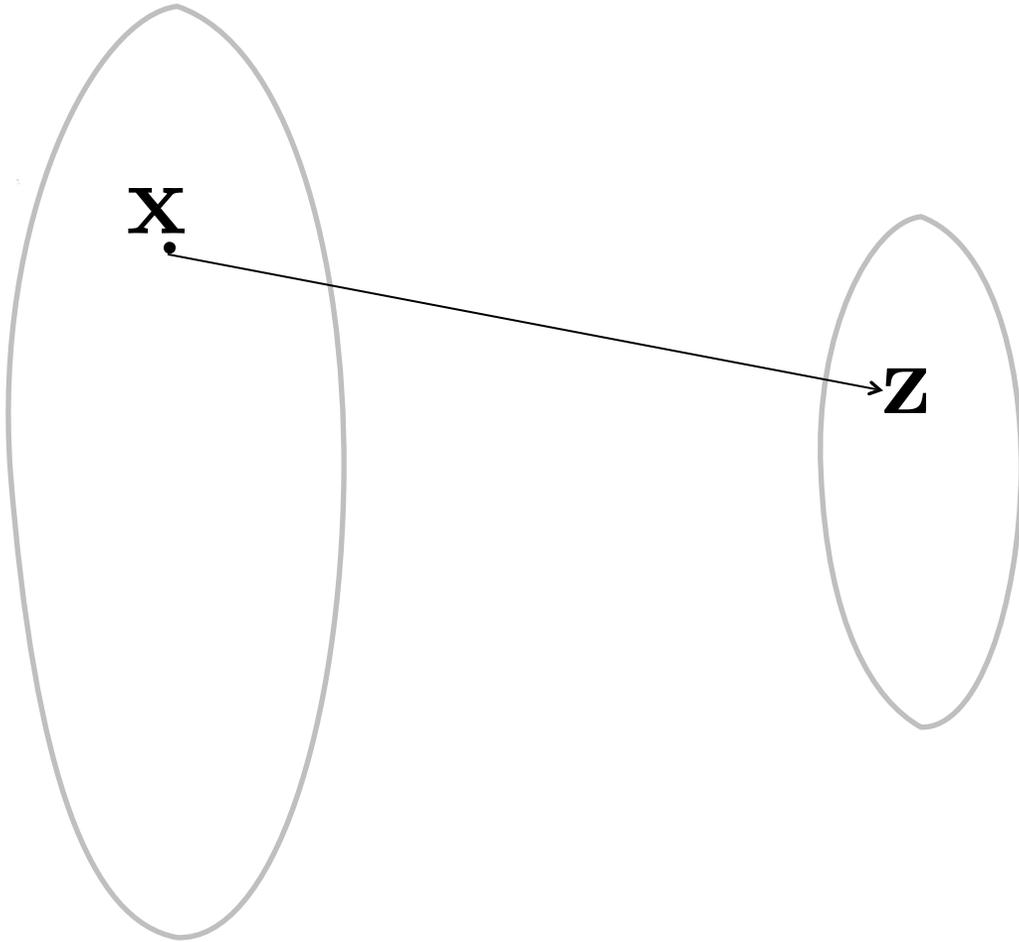
Autoregressive,
Normalizing
Flows,
Invertible NN,
and EBM

Diffusion
Models

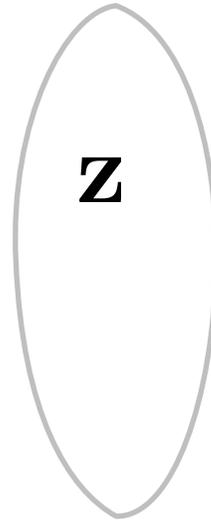
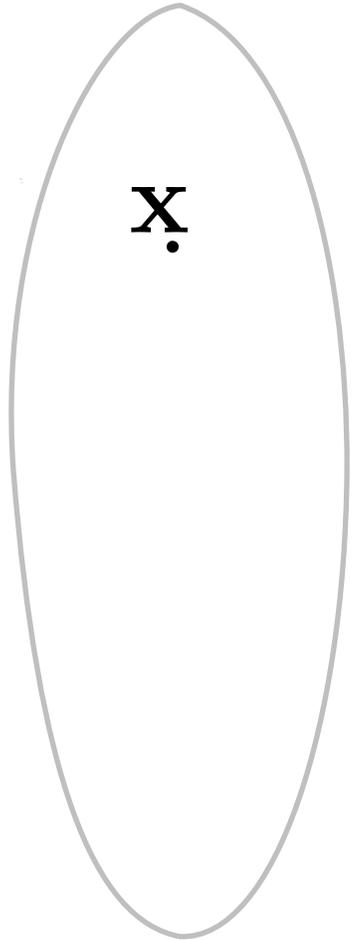
Robust Model $p(z|\mathbf{x})$



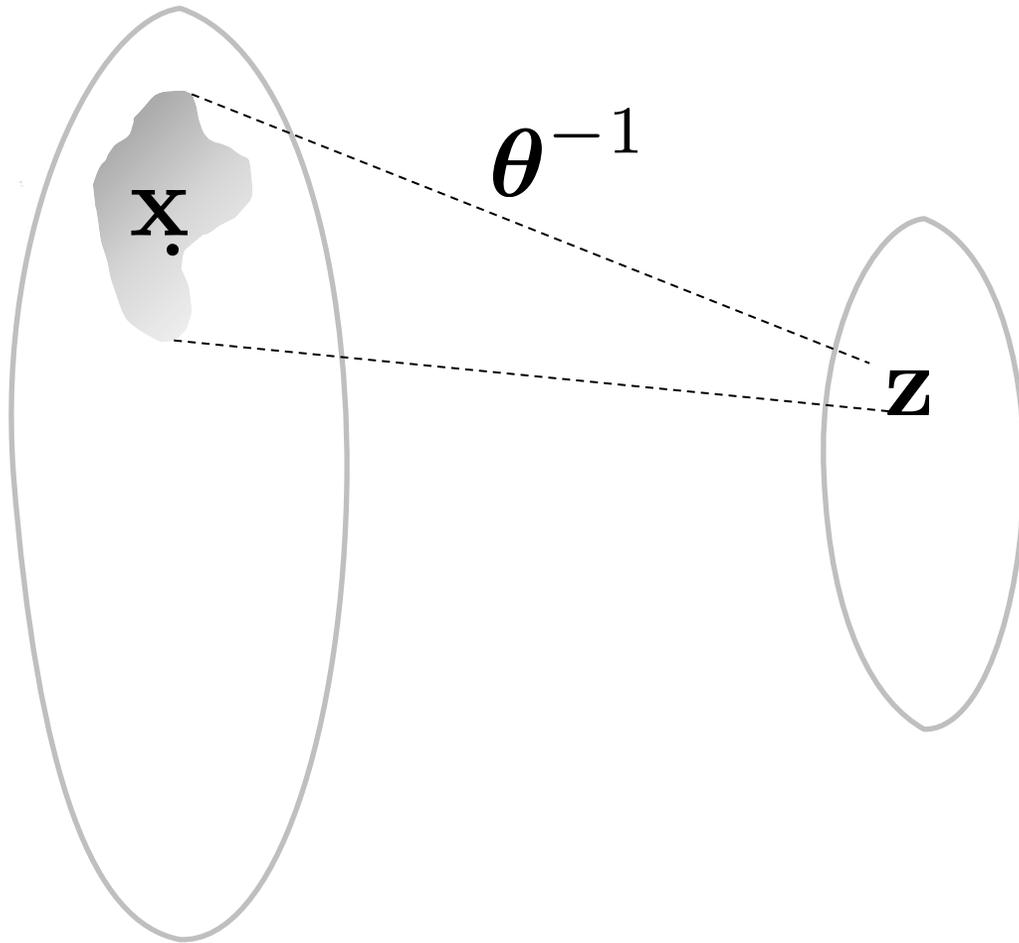
Robust Model $p(z|\mathbf{x})$



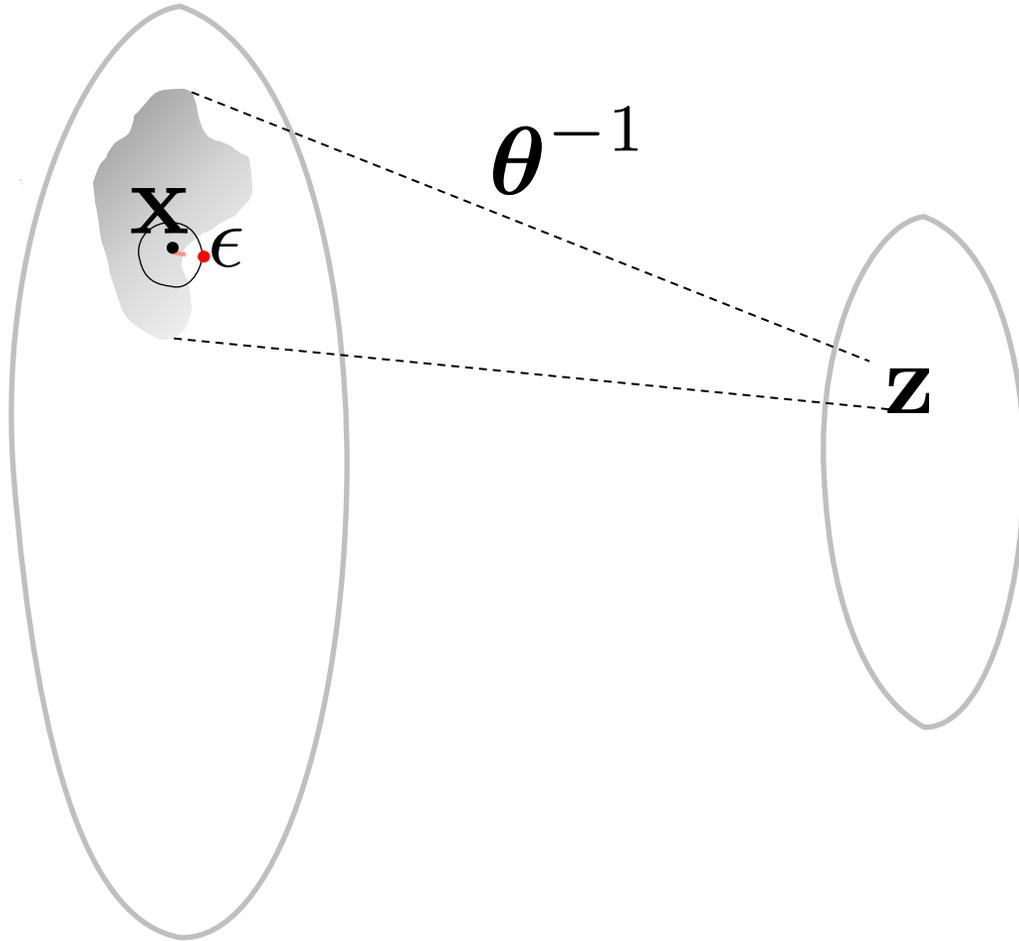
Robust Model $p(z|\mathbf{x})$



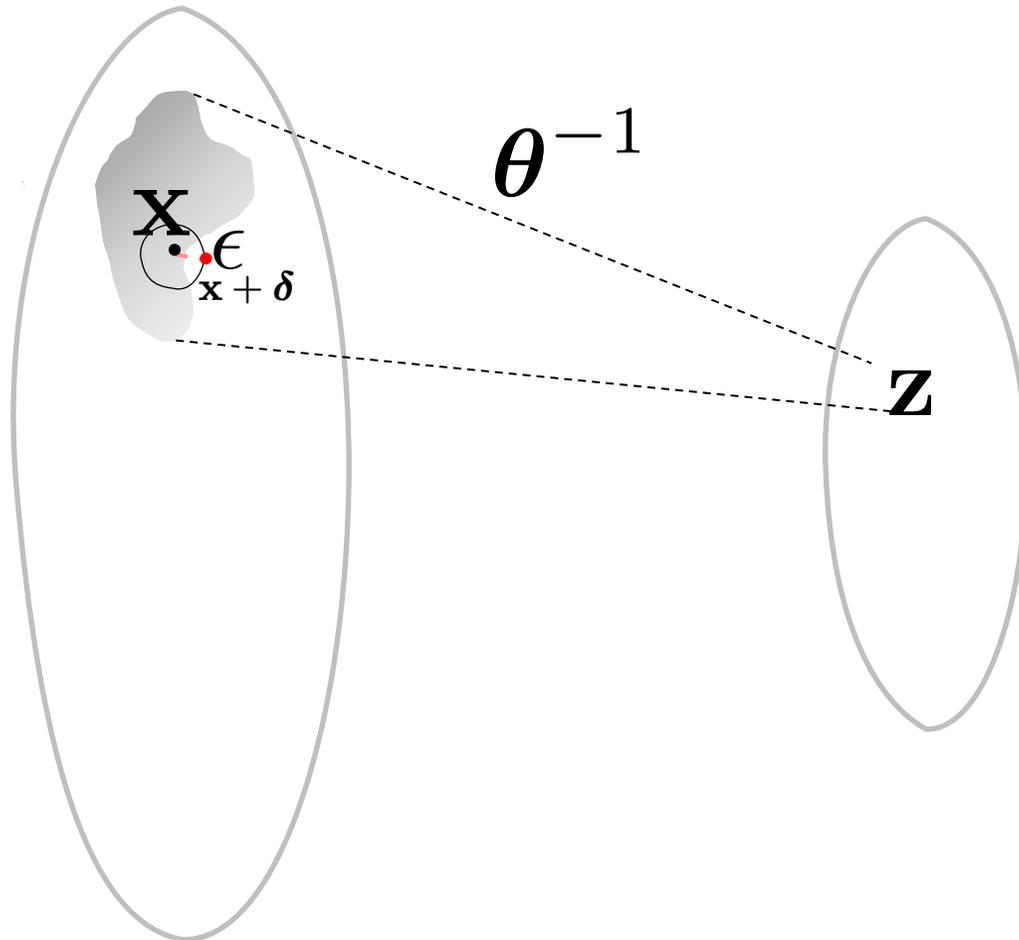
Robust Model $p(z|\mathbf{x})$



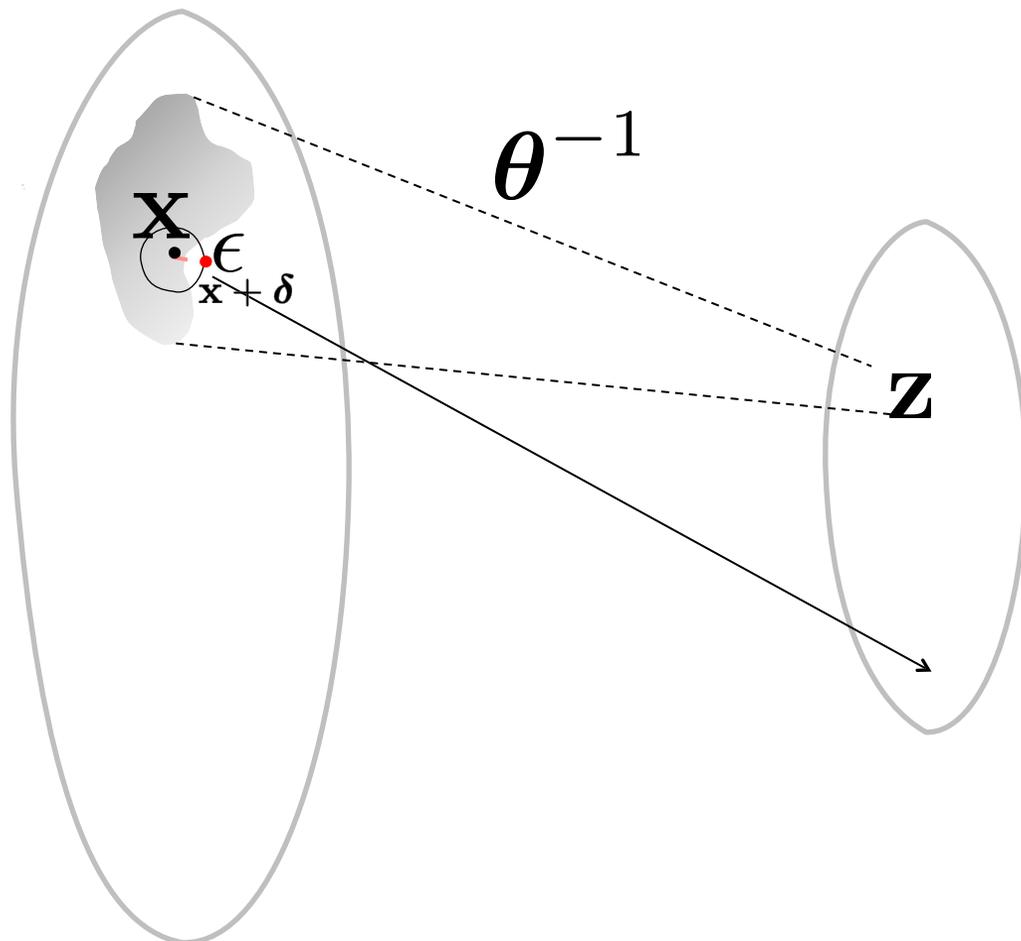
Robust Model $p(z|\mathbf{x})$



Robust Model $p(z|\mathbf{x})$

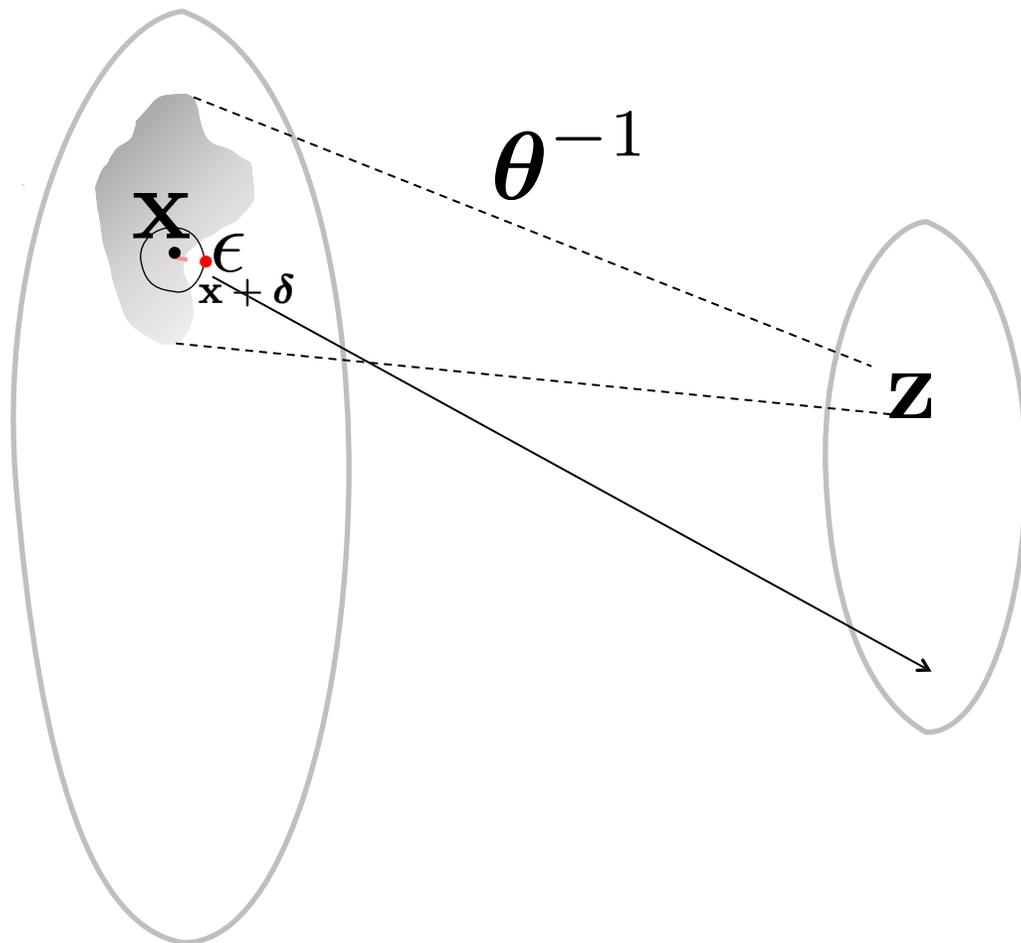


Robust Model $p(z|\mathbf{x})$



$$\delta^* = \arg \max_{\|\delta\|_p < \epsilon} \ell(\theta(\mathbf{x} + \delta), y)$$

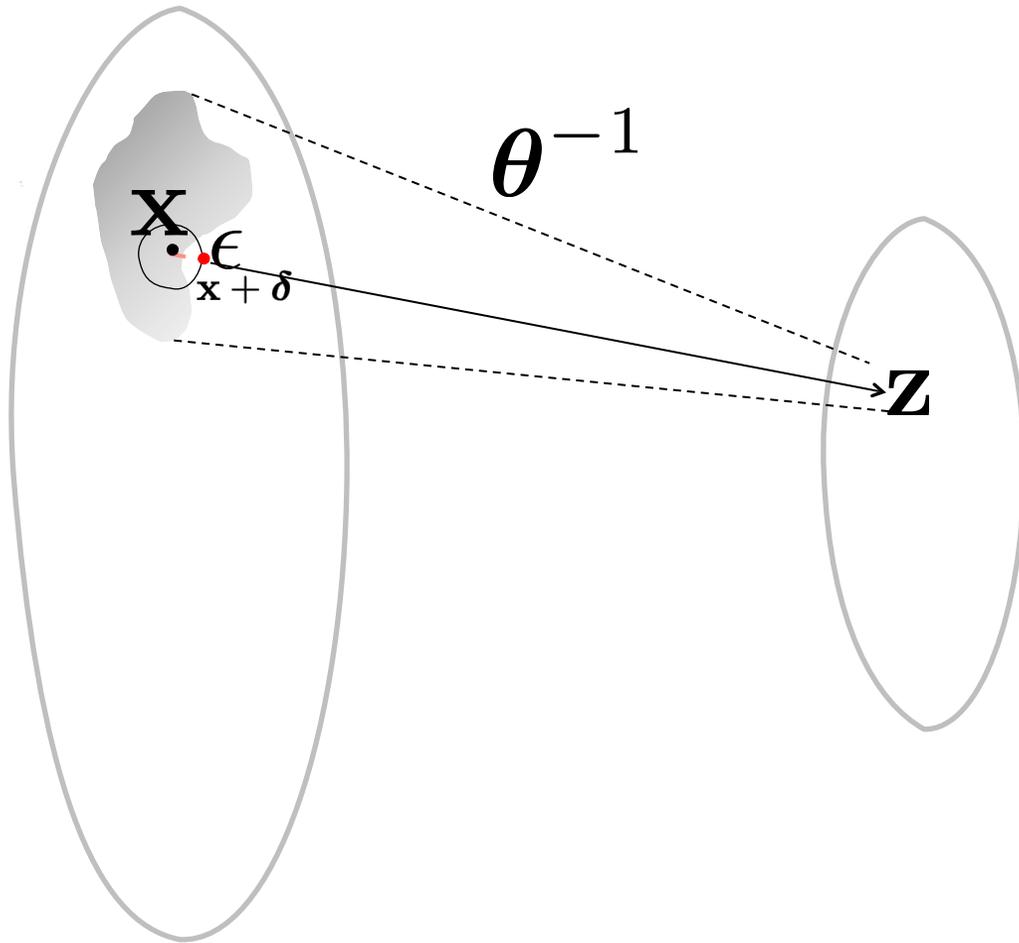
Robust Model $p(z|\mathbf{x})$



$$\theta^* = \arg \min_{\theta} \ell(\theta(\mathbf{x} + \delta^*), y) \quad \text{where}$$

$$\delta^* = \arg \max_{\|\delta\|_p < \epsilon} \ell(\theta(\mathbf{x} + \delta), y)$$

Robust Model $p(z|\mathbf{x})$

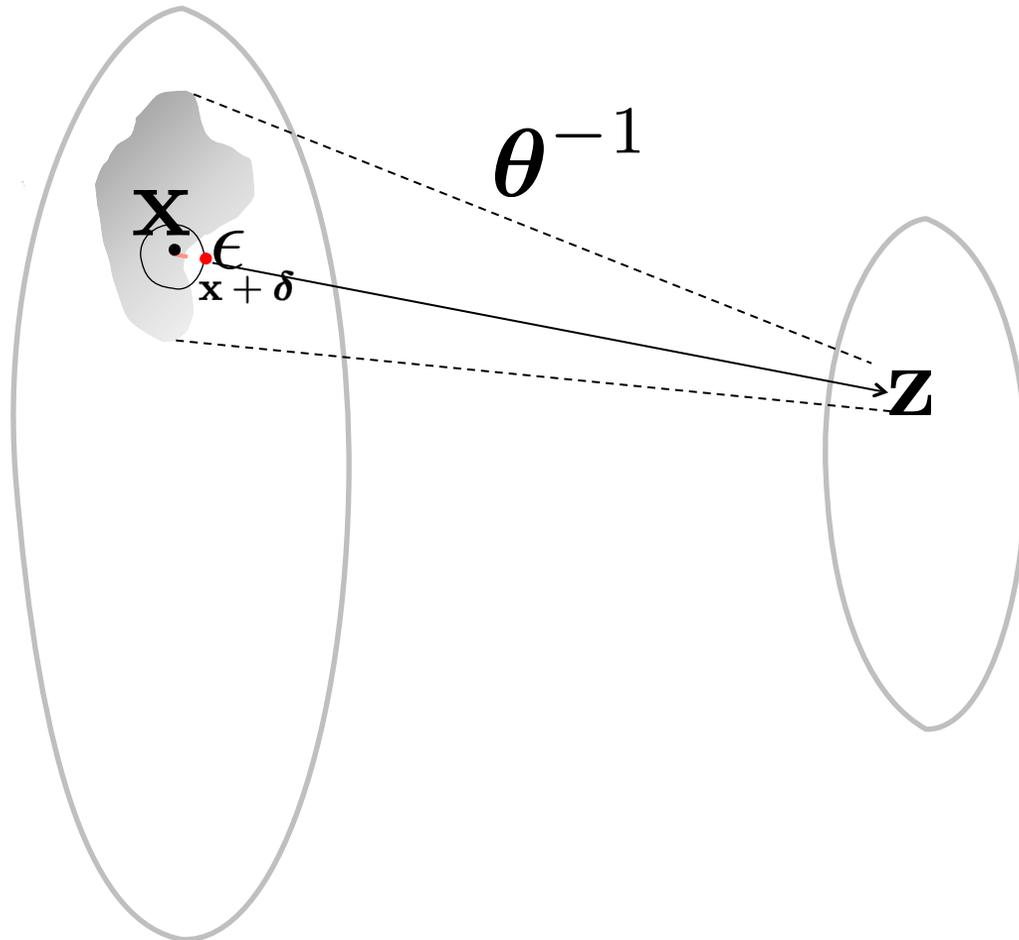


“adversarial training”

$$\theta^* = \arg \min_{\theta} \ell(\theta(\mathbf{x} + \delta^*), y) \quad \text{where}$$

$$\delta^* = \arg \max_{\|\delta\|_p < \epsilon} \ell(\theta(\mathbf{x} + \delta), y)$$

Robust Model $p(z|\mathbf{x})$



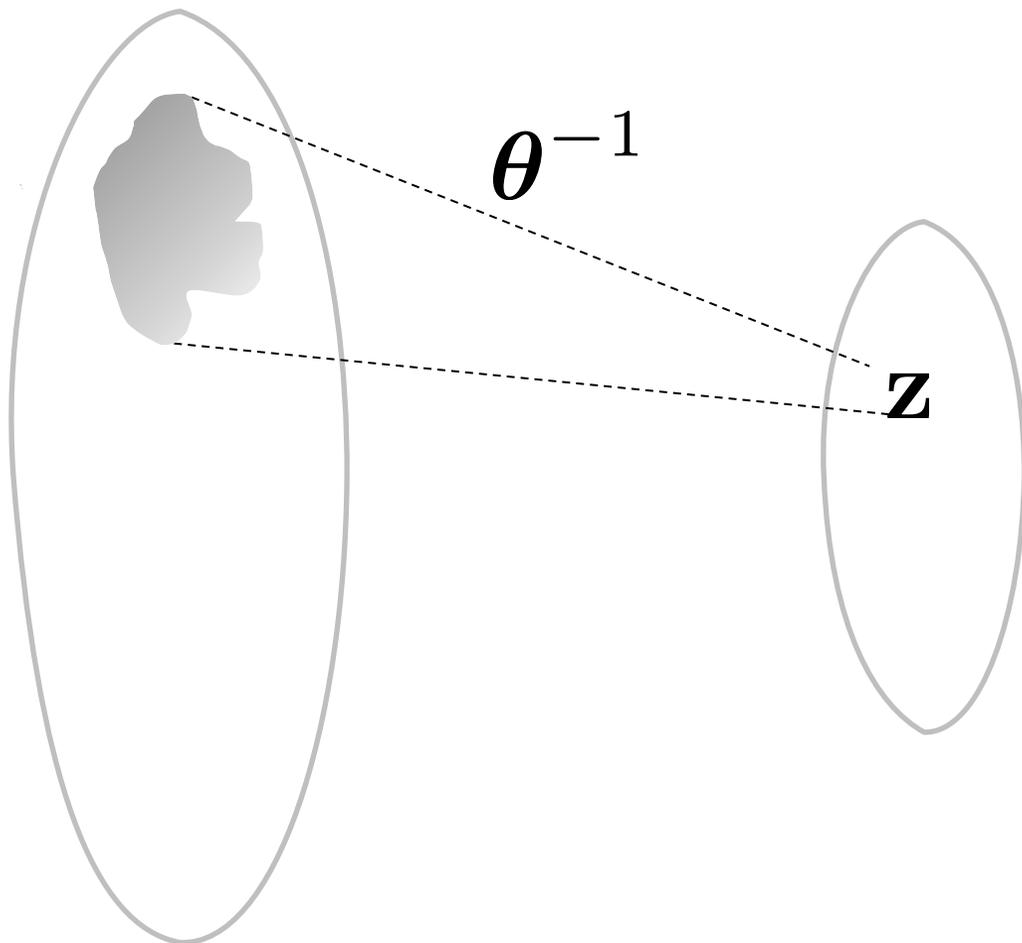
“adversarial training”

$$\theta^* = \arg \min_{\theta} \ell(\theta(\mathbf{x} + \delta^*), y) \quad \text{where}$$

$$\delta^* = \arg \max_{\|\delta\|_p < \epsilon} \ell(\theta(\mathbf{x} + \delta), y)$$

⊖ Decreases the accuracy on natural data

Robust Model $p(z|\mathbf{x})$



“adversarial training”

$$\theta^* = \arg \min_{\theta} \ell(\theta(\mathbf{x} + \delta^*), y) \quad \text{where}$$

$$\delta^* = \arg \max_{\|\delta\|_p < \epsilon} \ell(\theta(\mathbf{x} + \delta), y)$$

- ⊖ Decreases the accuracy on natural data
- ⊕ Develops “generative” behavior

Looking at the input gradients



Looking at the input gradients



$$\nabla_{\mathbf{x}} \ell_{\text{CE}}(\mathbf{x}, z; \boldsymbol{\theta})$$

Looking at the input gradients

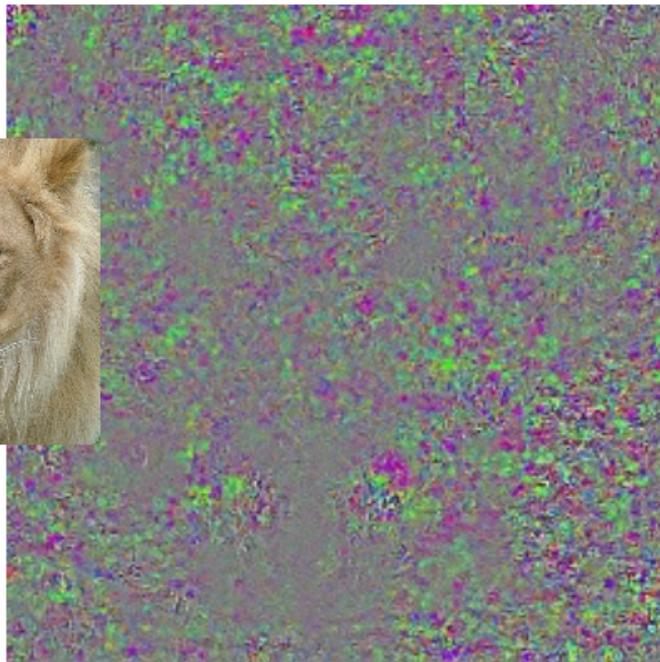
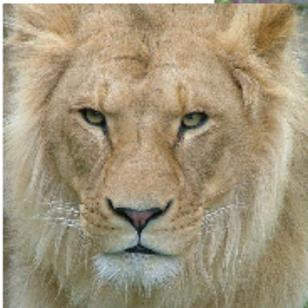


$$\nabla_{\mathbf{x}} \ell_{\text{CE}}(\mathbf{x}, z; \boldsymbol{\theta})$$

Standard, non-robust

Wang et al. [4]

Input



Looking at the input gradients

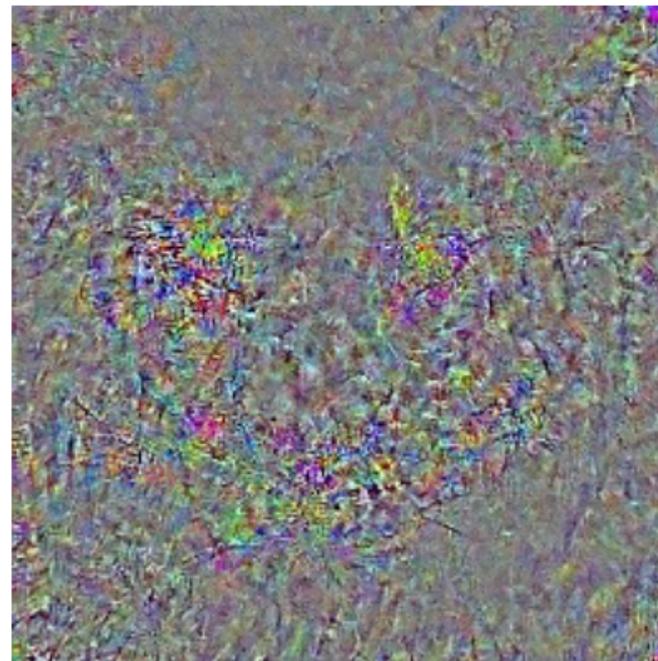
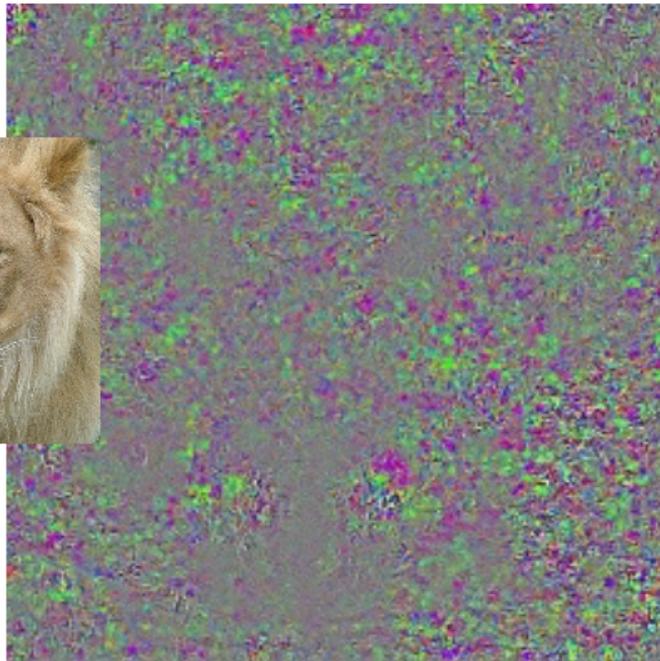
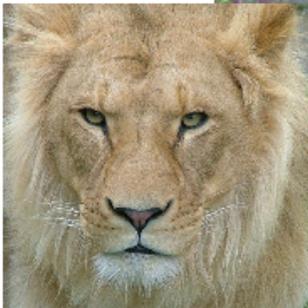
$$\nabla_{\mathbf{x}} \ell_{\text{CE}}(\mathbf{x}, z; \boldsymbol{\theta})$$

Standard, non-robust

Wang et al. [4]

$\ell_2, \epsilon=0.01$

Input



Looking at the input gradients

$$\nabla_{\mathbf{x}} \ell_{\text{CE}}(\mathbf{x}, z; \boldsymbol{\theta})$$

Standard, non-robust

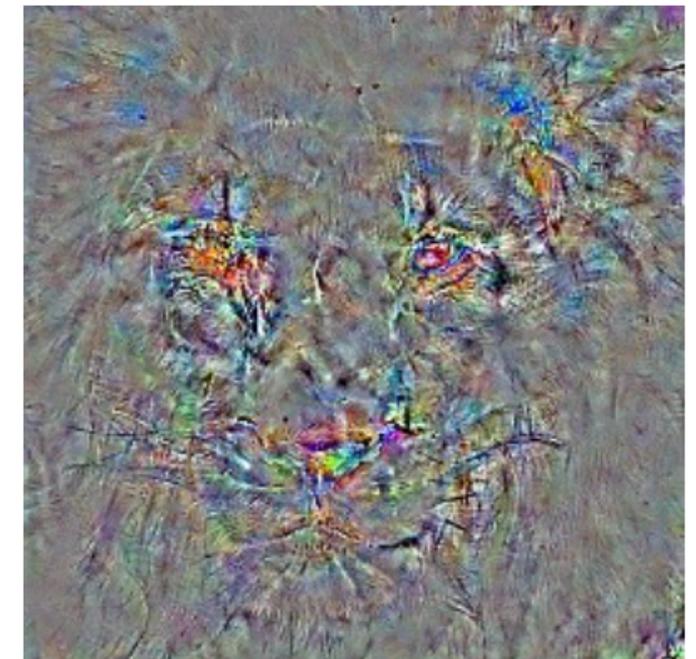
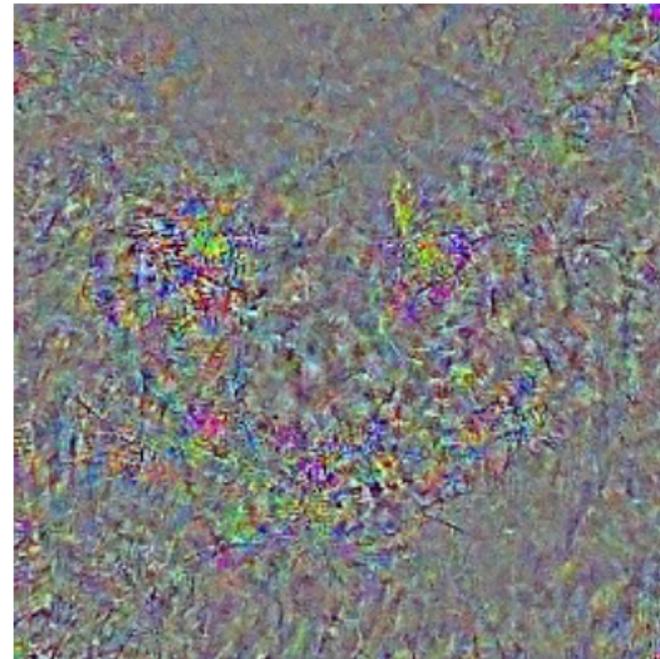
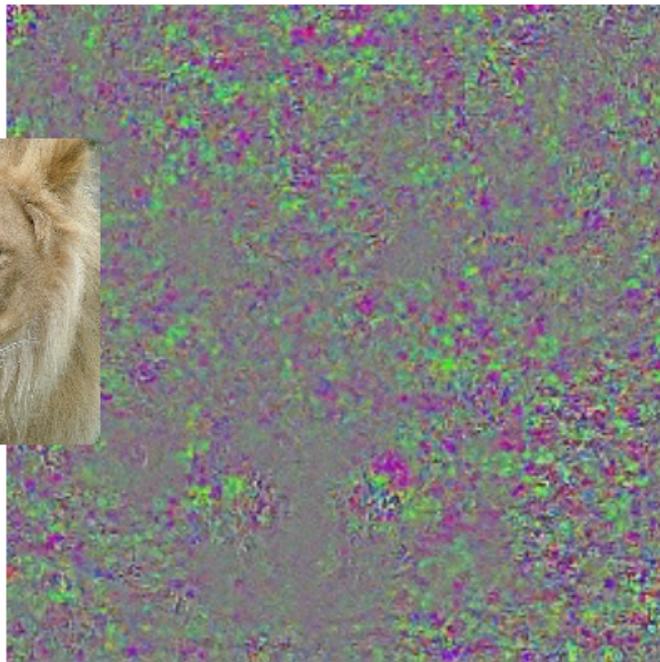
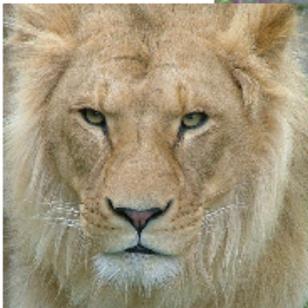
“Robust” family

Wang et al. [4]

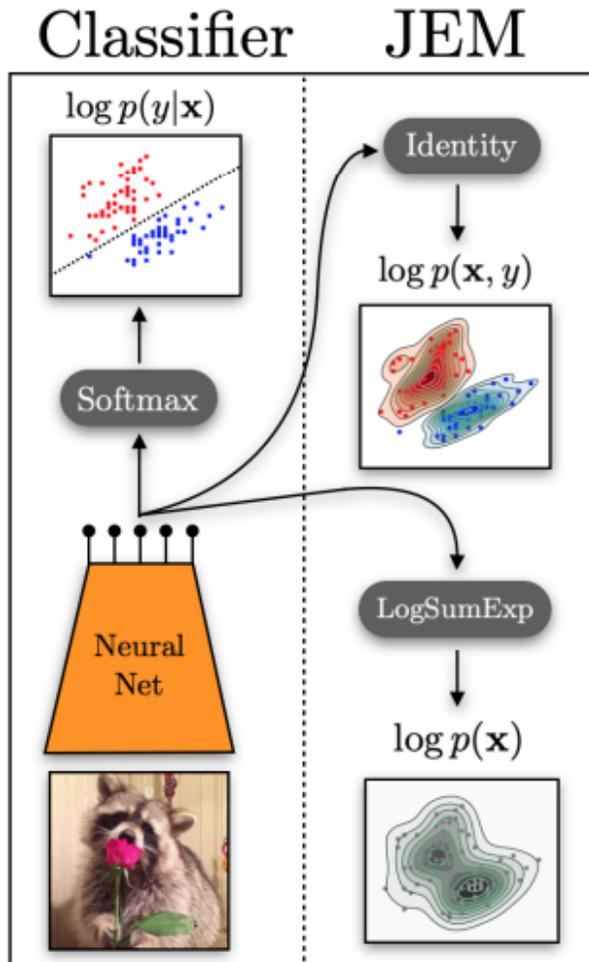
$\ell_2, \epsilon=0.01$

$\ell_2, \epsilon=0.05$

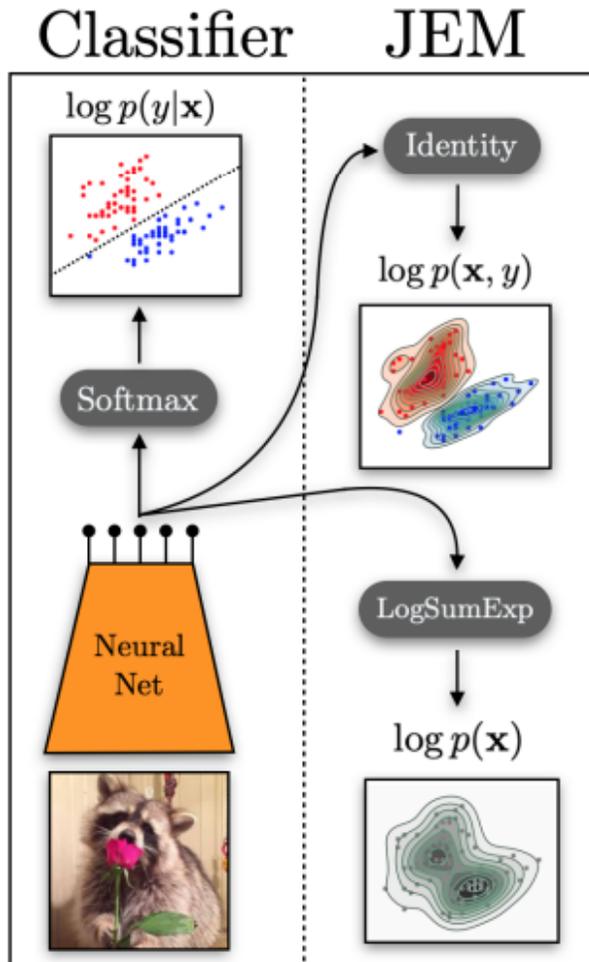
Input



Why Robust Models behave as Generative?



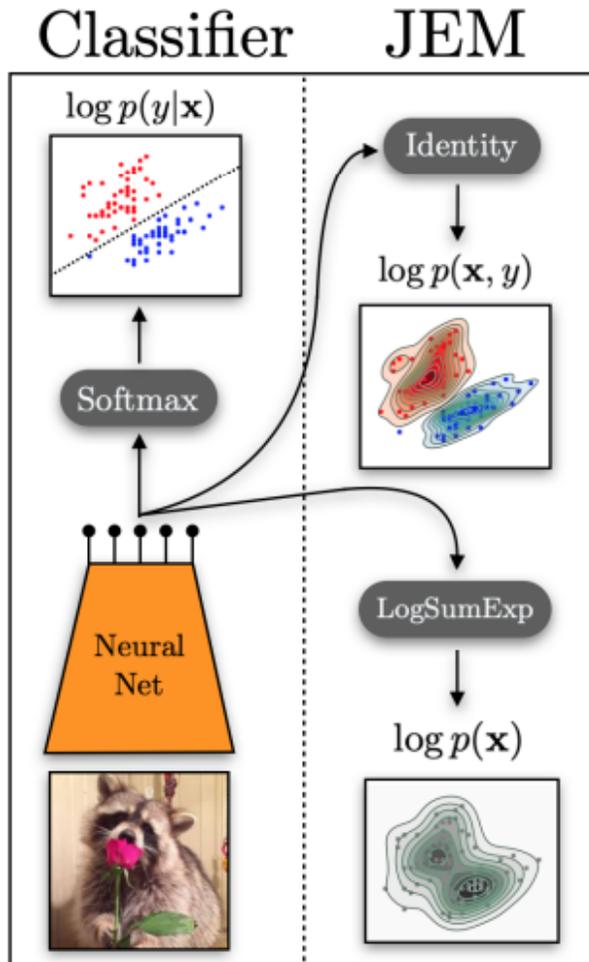
[image is from JEM - ICLR20]



[image is from JEM - ICLR20]

“energy-based model”

$$p_{\theta}(\mathbf{x}) = \frac{\exp(-E_{\theta}(\mathbf{x}))}{P_{\theta}}$$



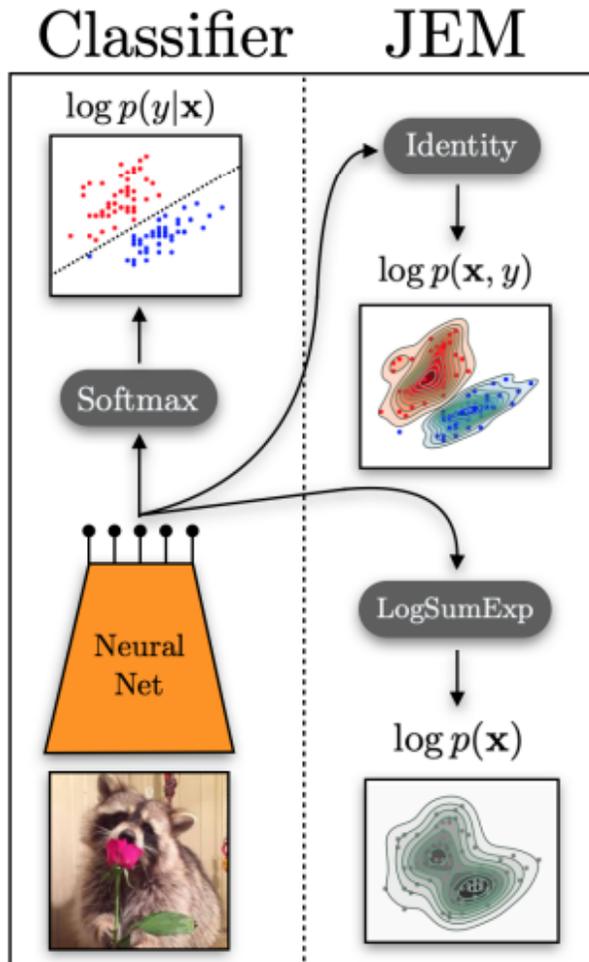
[image is from JEM - ICLR20]

“energy-based model”

$$p_{\theta}(\mathbf{x}) = \frac{\exp(-E_{\theta}(\mathbf{x}))}{P_{\theta}}$$

Classifier

$$p(z = i|\mathbf{x}) = \frac{\exp F_{\theta}(x)[i]}{\sum_{i=1}^K \exp F_{\theta}(x)[i]}$$



[image is from JEM - ICLR20]

“energy-based model”

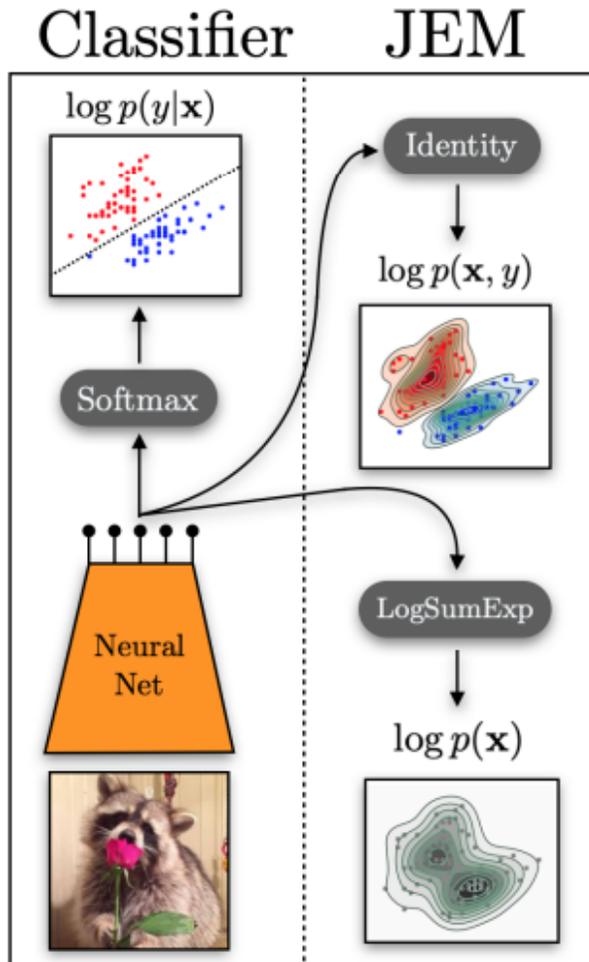
$$p_{\theta}(\mathbf{x}) = \frac{\exp(-E_{\theta}(\mathbf{x}))}{P_{\theta}}$$

Classifier

$$p(z = i|\mathbf{x}) = \frac{\exp F_{\theta}(x)[i]}{\sum_{i=1}^K \exp F_{\theta}(x)[i]}$$

$$E_{\theta}(\mathbf{x}, z)$$

Joint energy: datum vs label



[image is from JEM - ICLR20]

“energy-based model”

$$p_{\theta}(\mathbf{x}) = \frac{\exp(-E_{\theta}(\mathbf{x}))}{P_{\theta}}$$

Classifier

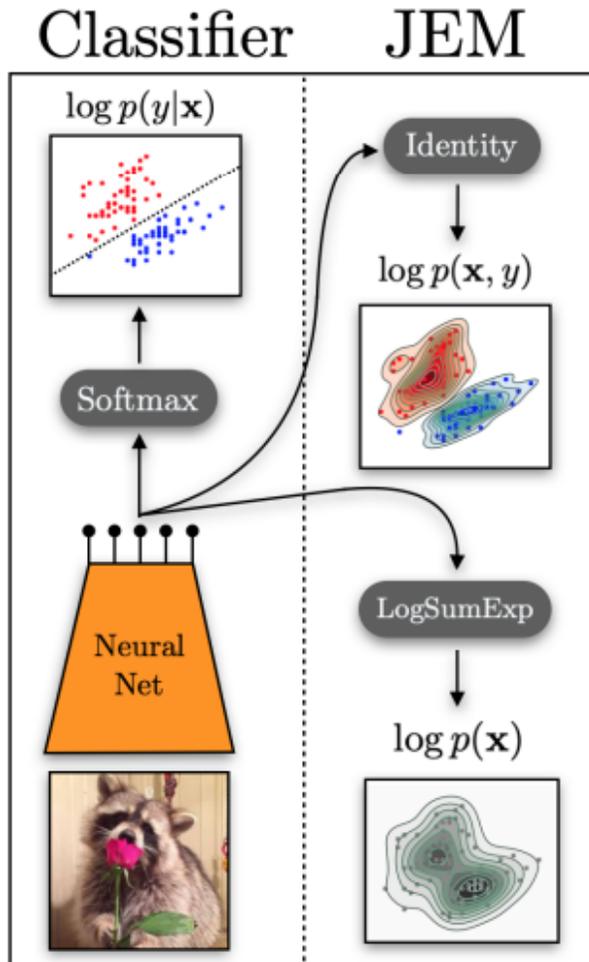
$$p(z = i | \mathbf{x}) = \frac{\exp F_{\theta}(x)[i]}{\sum_{i=1}^K \exp F_{\theta}(x)[i]}$$

$$E_{\theta}(\mathbf{x}, z)$$

Joint energy: datum vs label

$$E_{\theta}(\mathbf{x})$$

Marginal: Energy of datum



[image is from JEM - ICLR20]

“energy-based model”

$$p_{\theta}(\mathbf{x}) = \frac{\exp(-E_{\theta}(\mathbf{x}))}{P_{\theta}}$$

Classifier

$$p(z = i | \mathbf{x}) = \frac{\exp F_{\theta}(x)[i]}{\sum_{i=1}^K \exp F_{\theta}(x)[i]}$$

$$E_{\theta}(\mathbf{x}, z)$$

Joint energy: datum vs label

$$E_{\theta}(\mathbf{x})$$

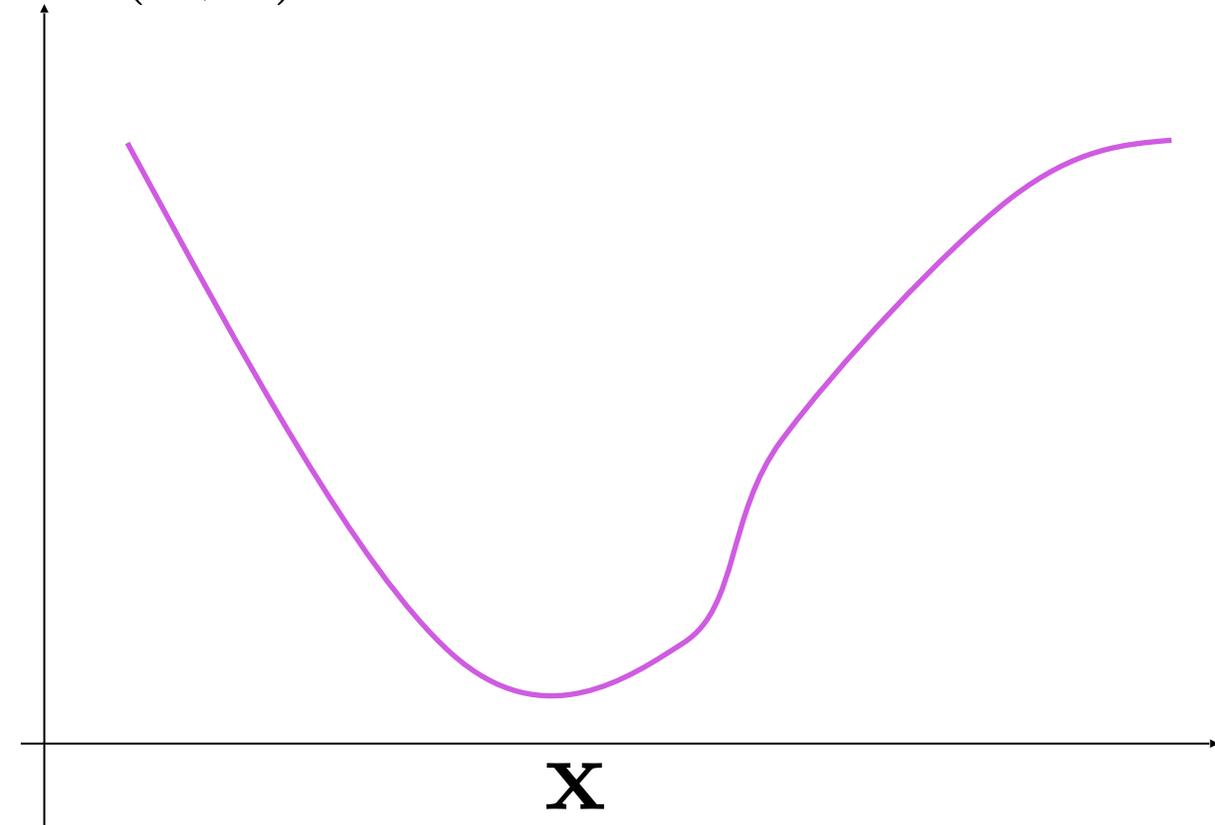
Marginal: Energy of datum

$$\ell_{\text{CE}}(\mathbf{x}, z; \theta) = E(\mathbf{x}, z; \theta) - E(\mathbf{x}; \theta)$$

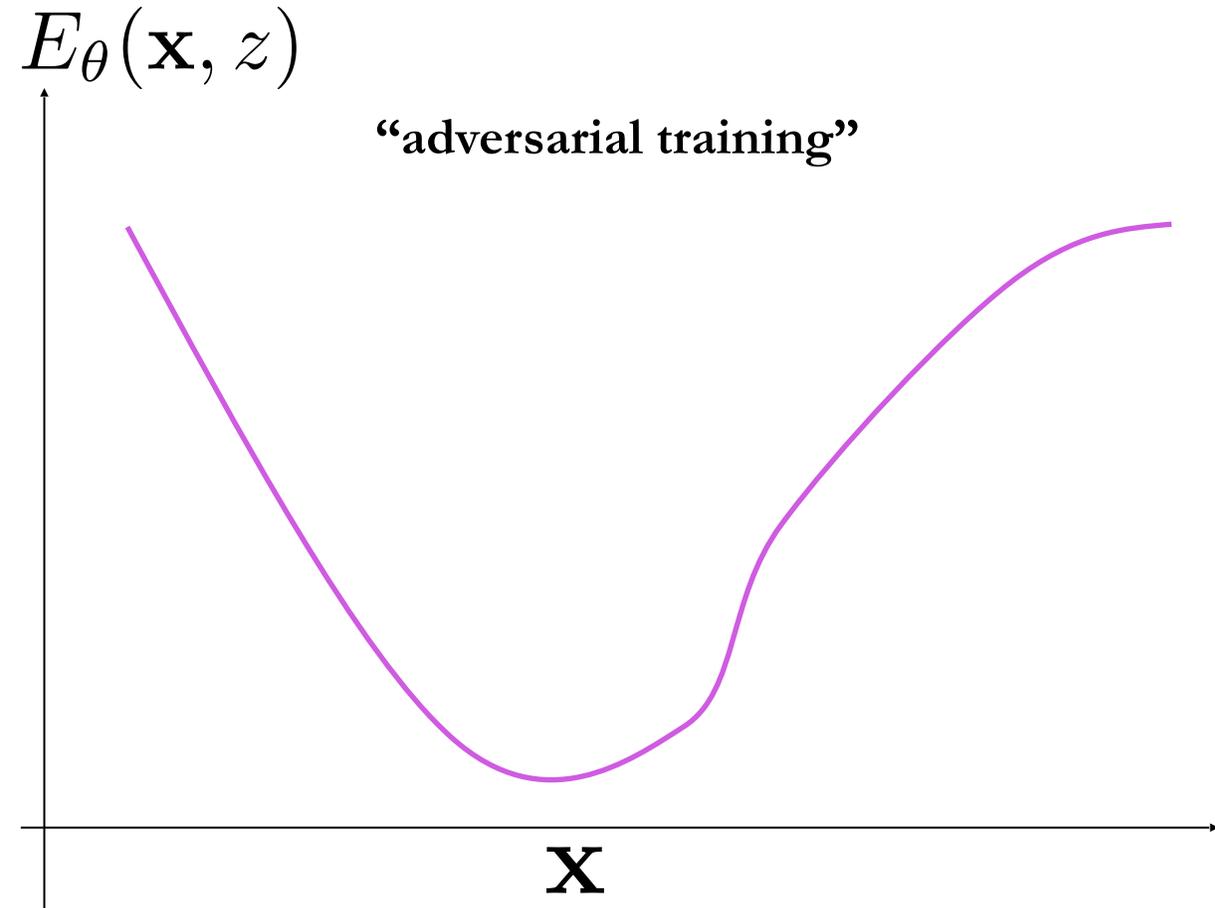
Why Robust Models behave as Generatives?

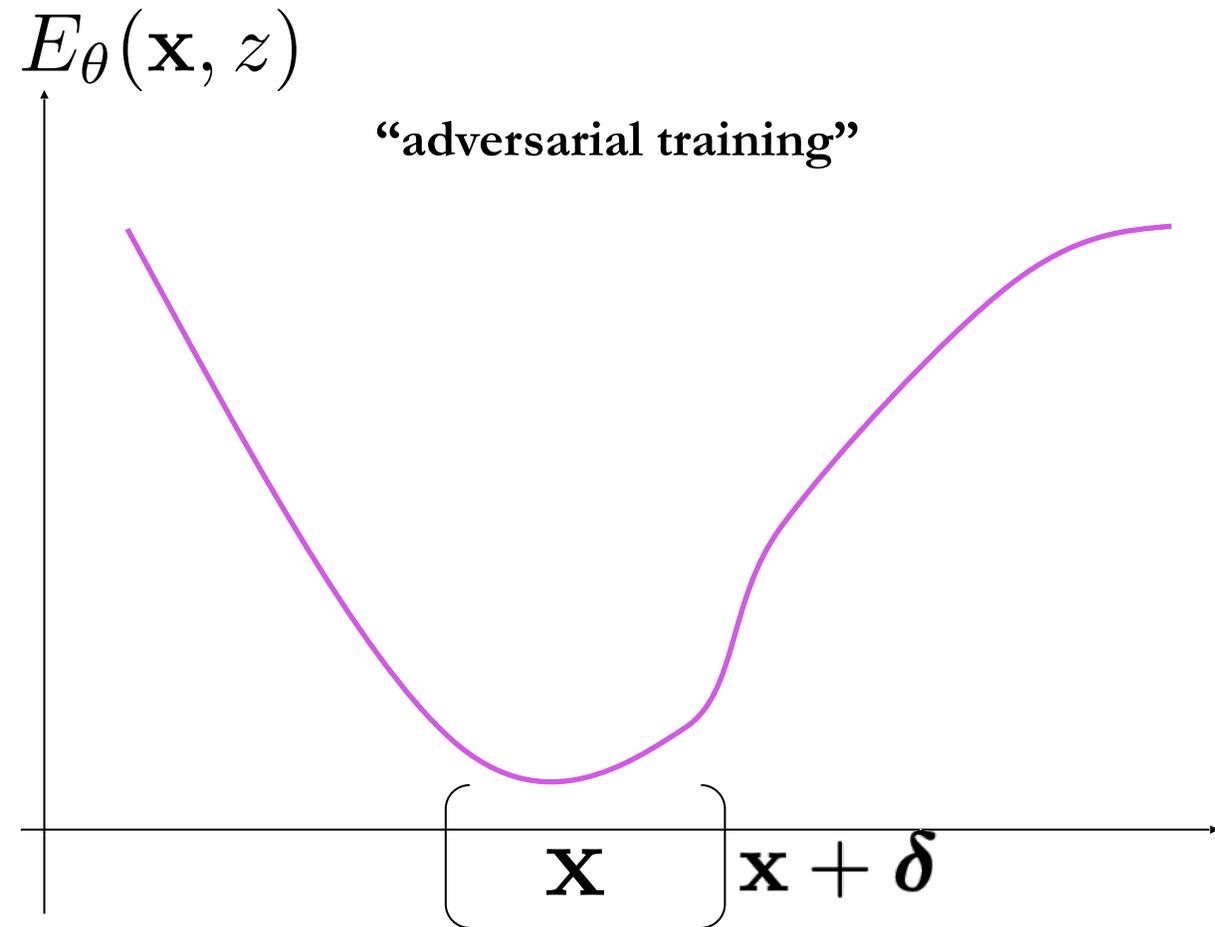


$$E_{\theta}(\mathbf{x}, z)$$

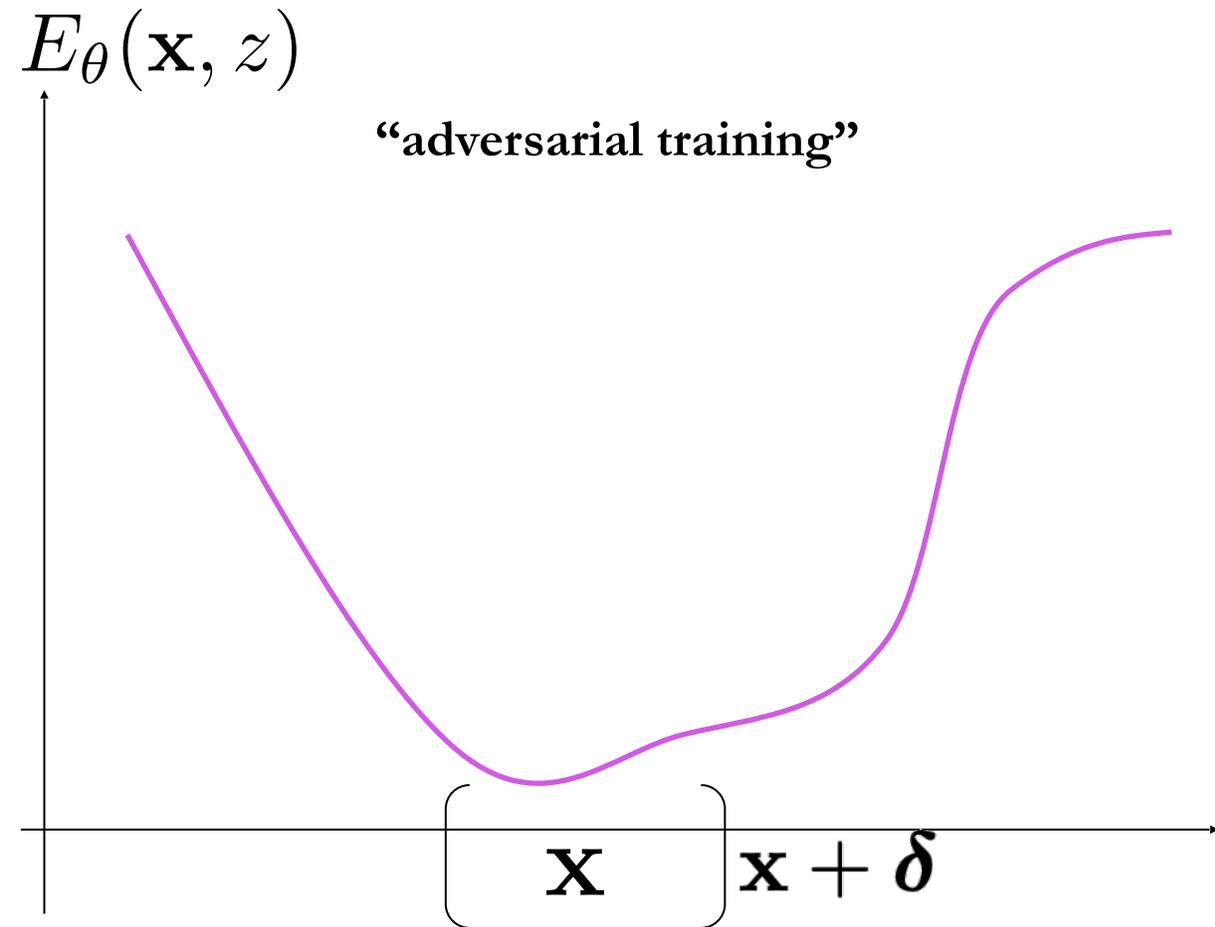


$$E_{\theta}(\mathbf{x}, z)$$



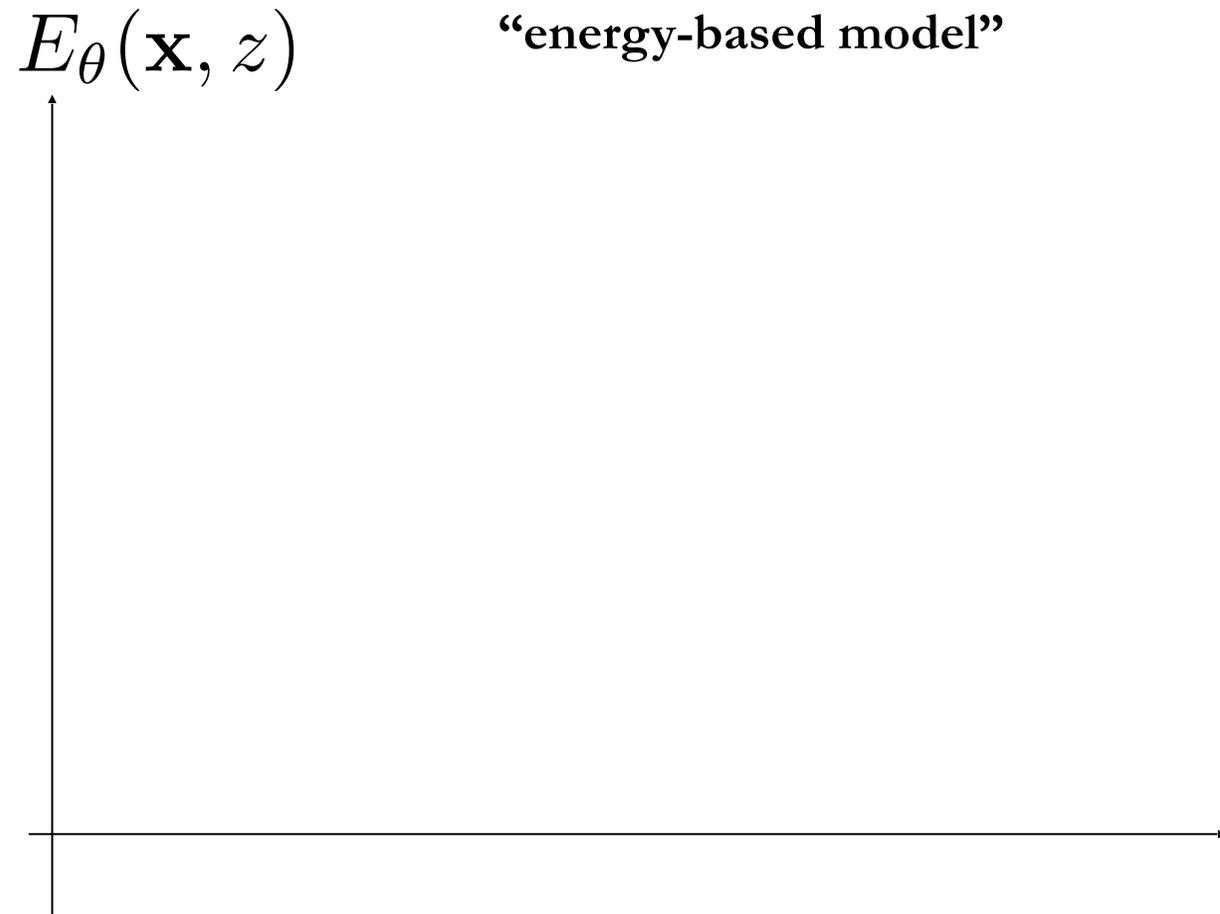
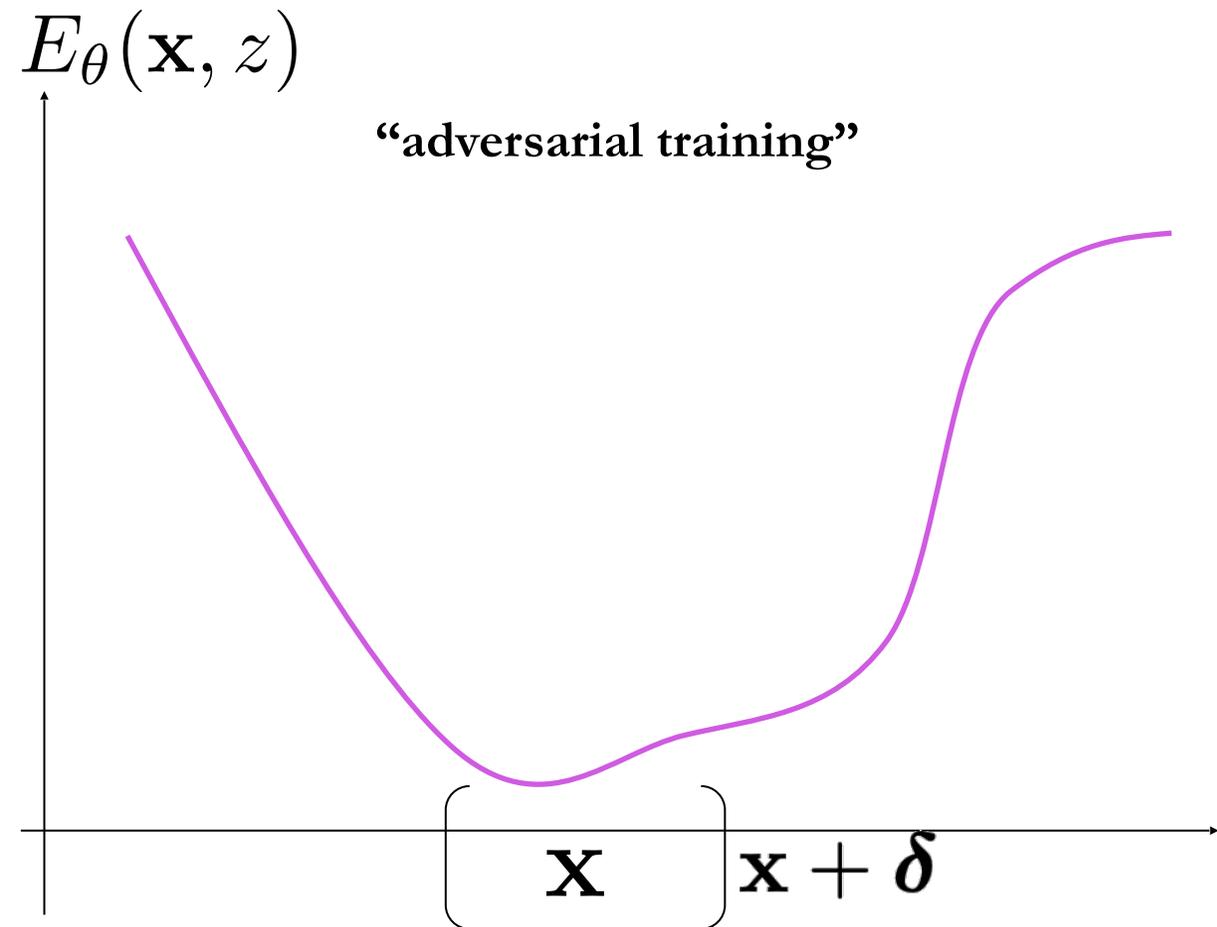


$E_{\theta}(\mathbf{x}, z)$

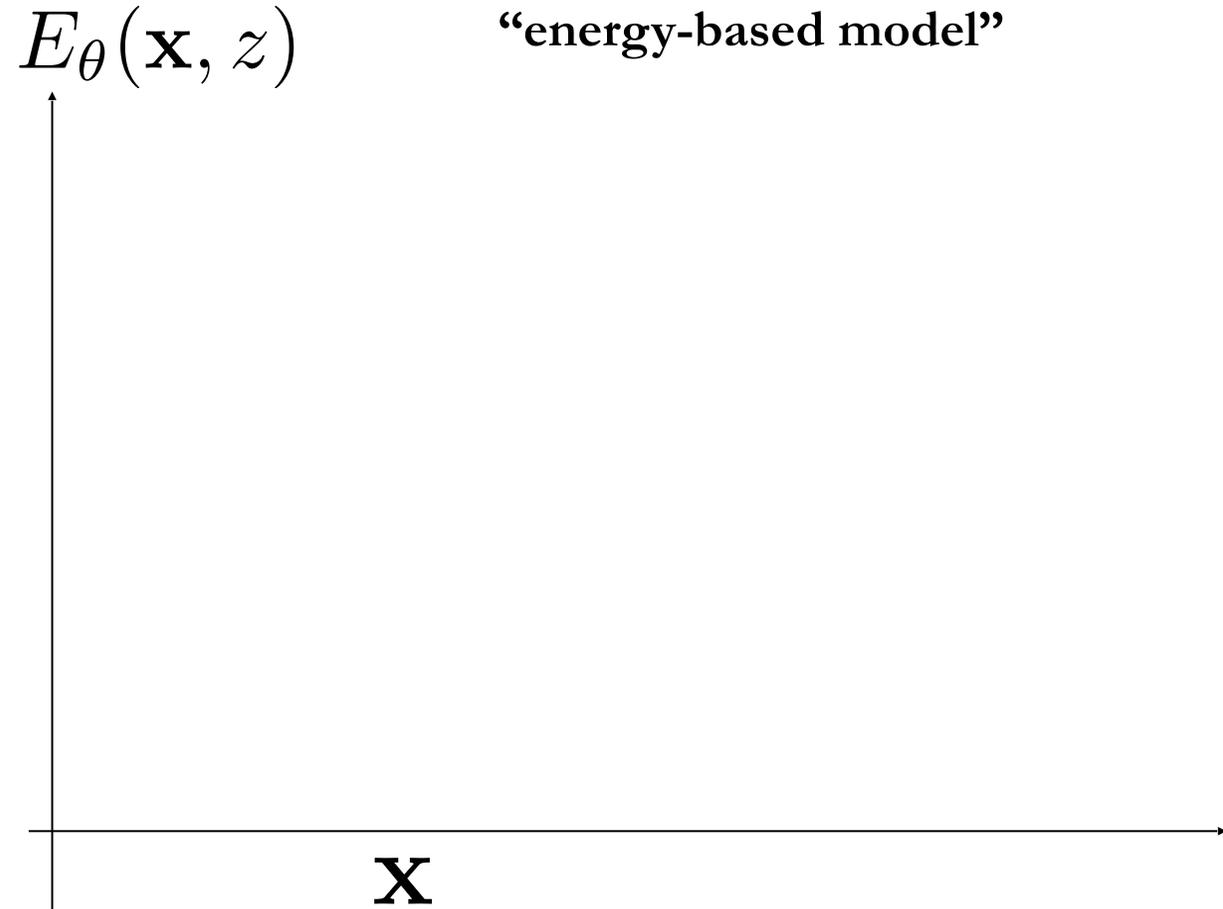
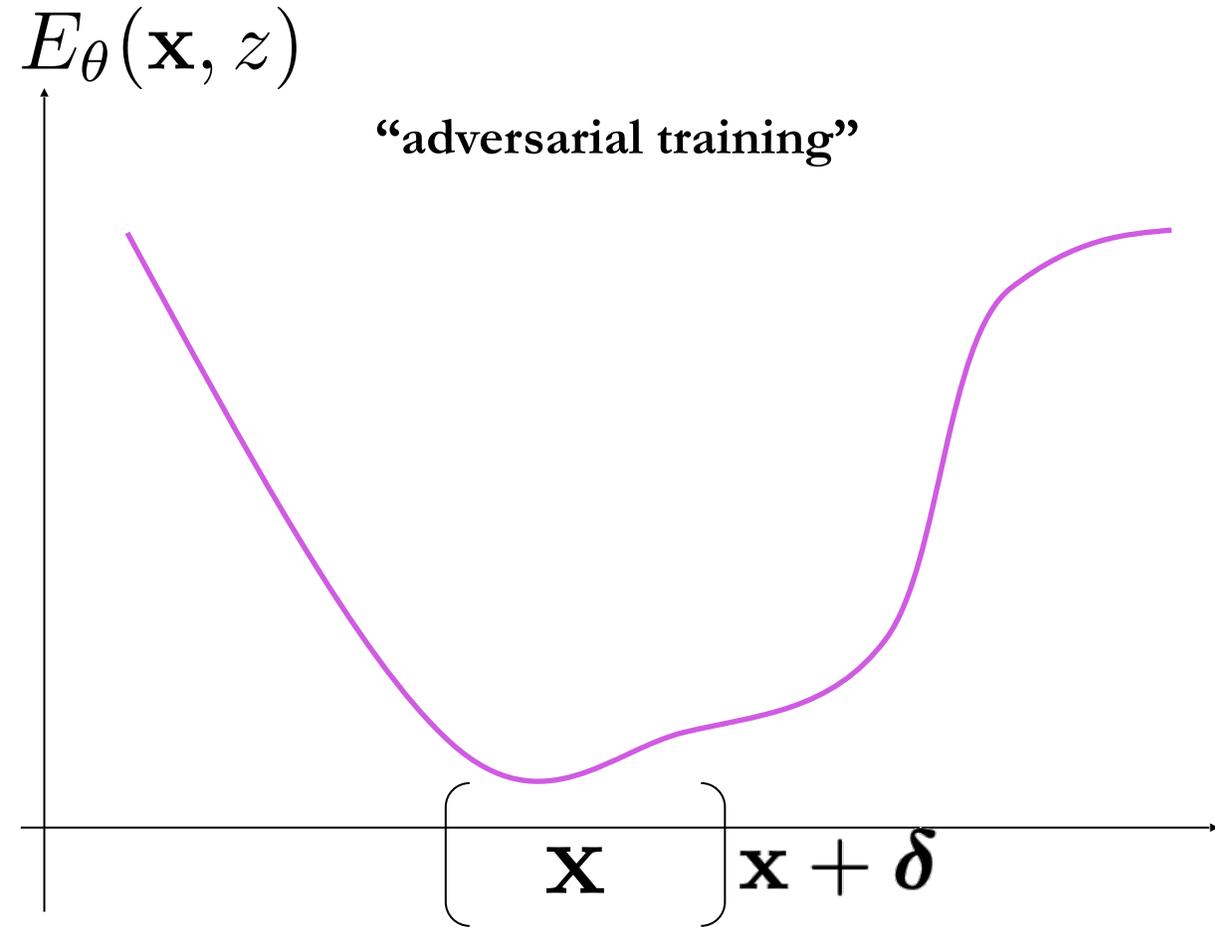


$E_{\theta}(\mathbf{x}, z)$

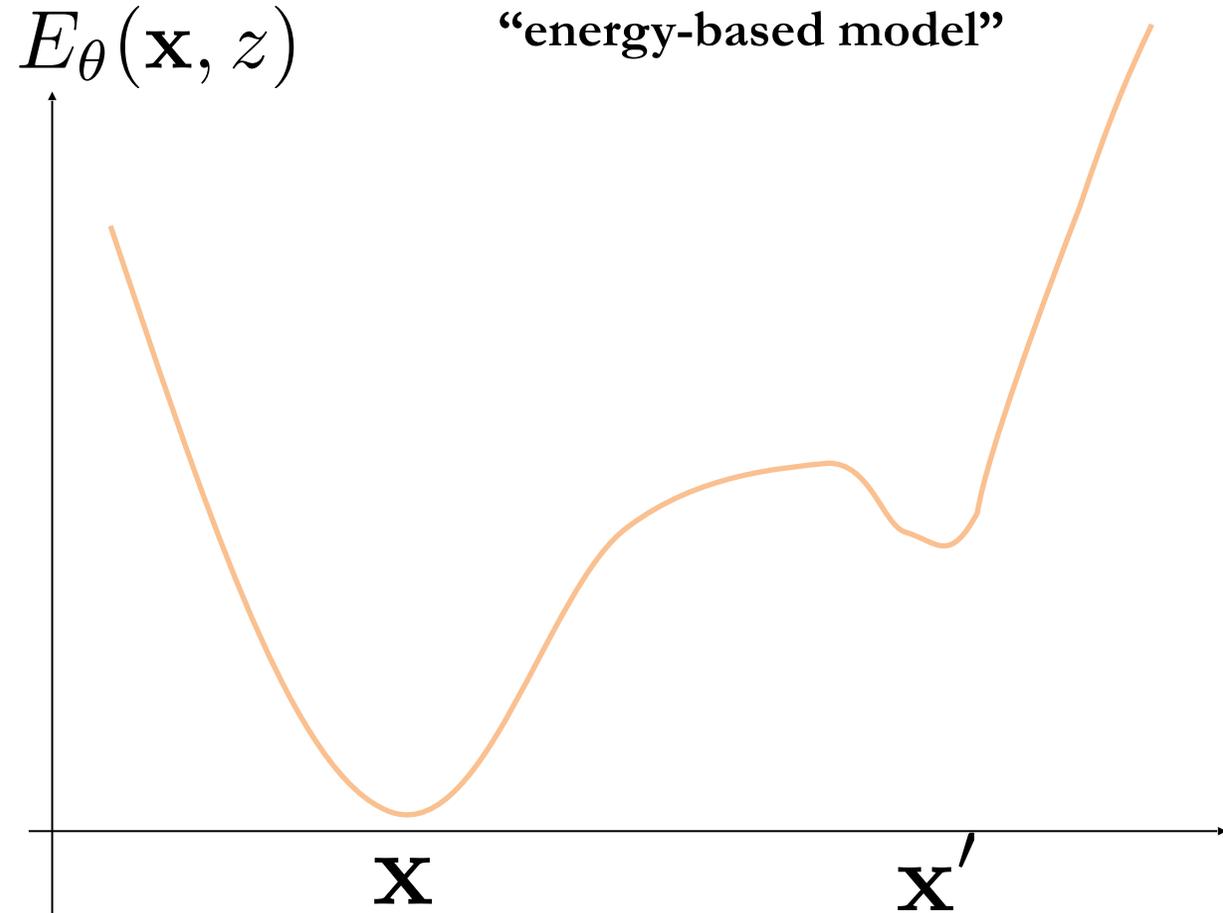
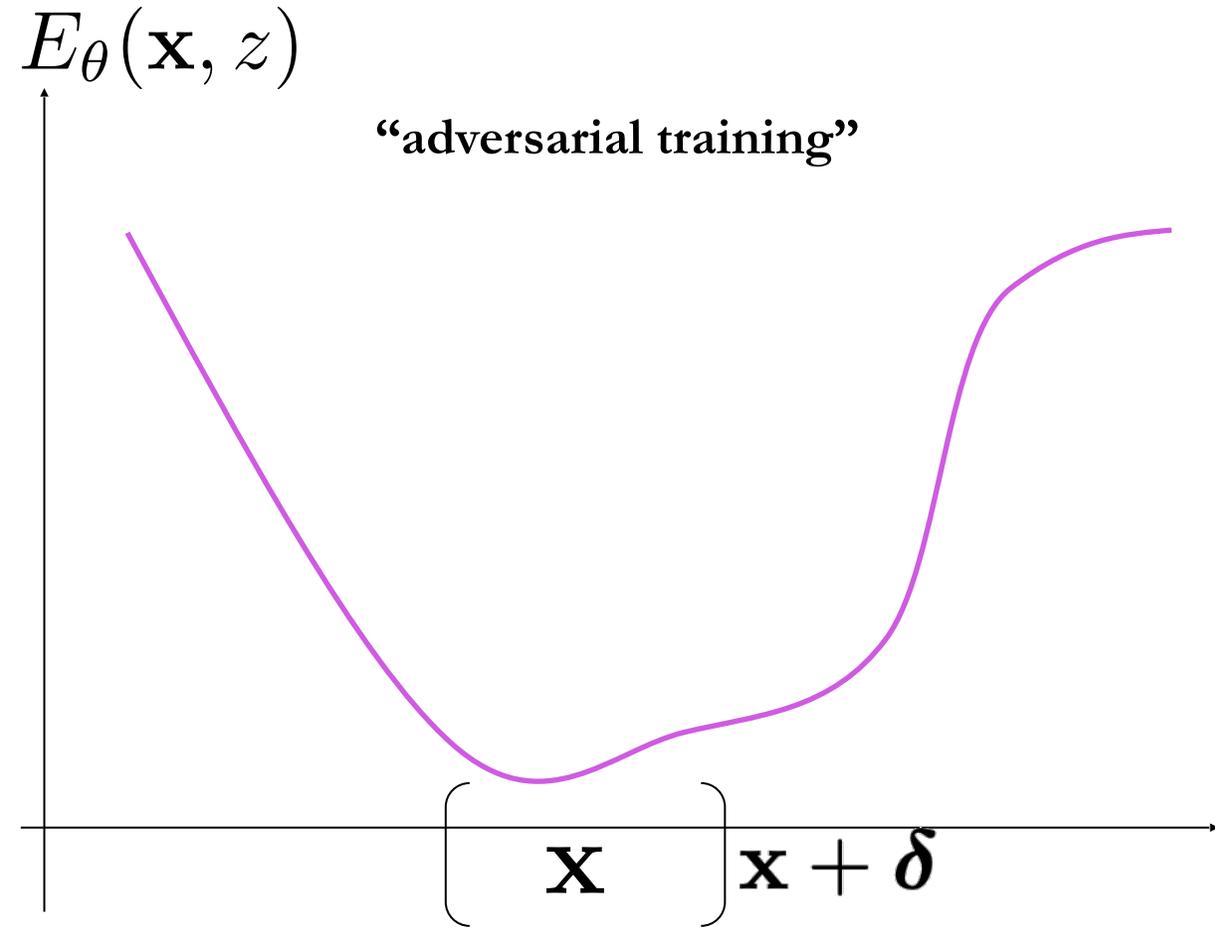
Why Robust Models behave as Generatives?

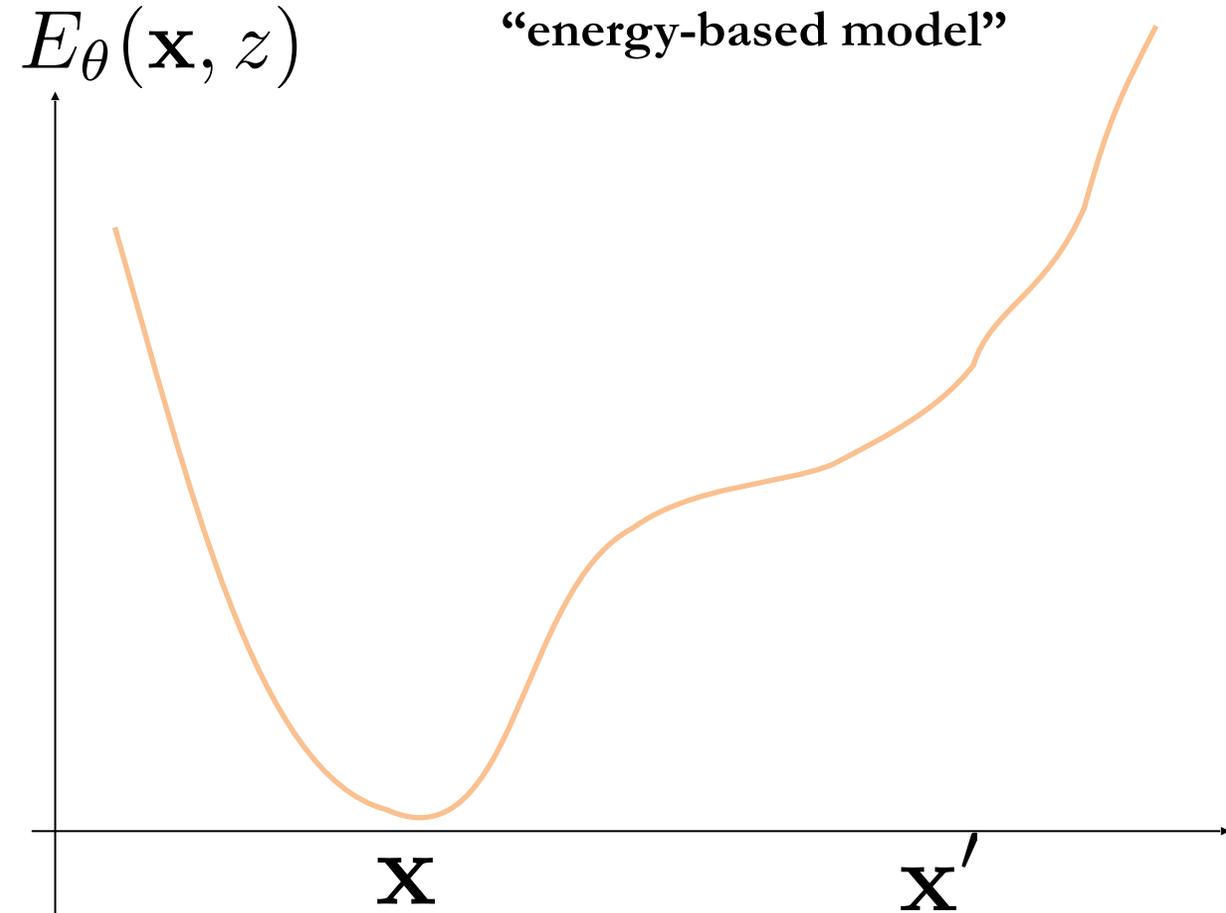
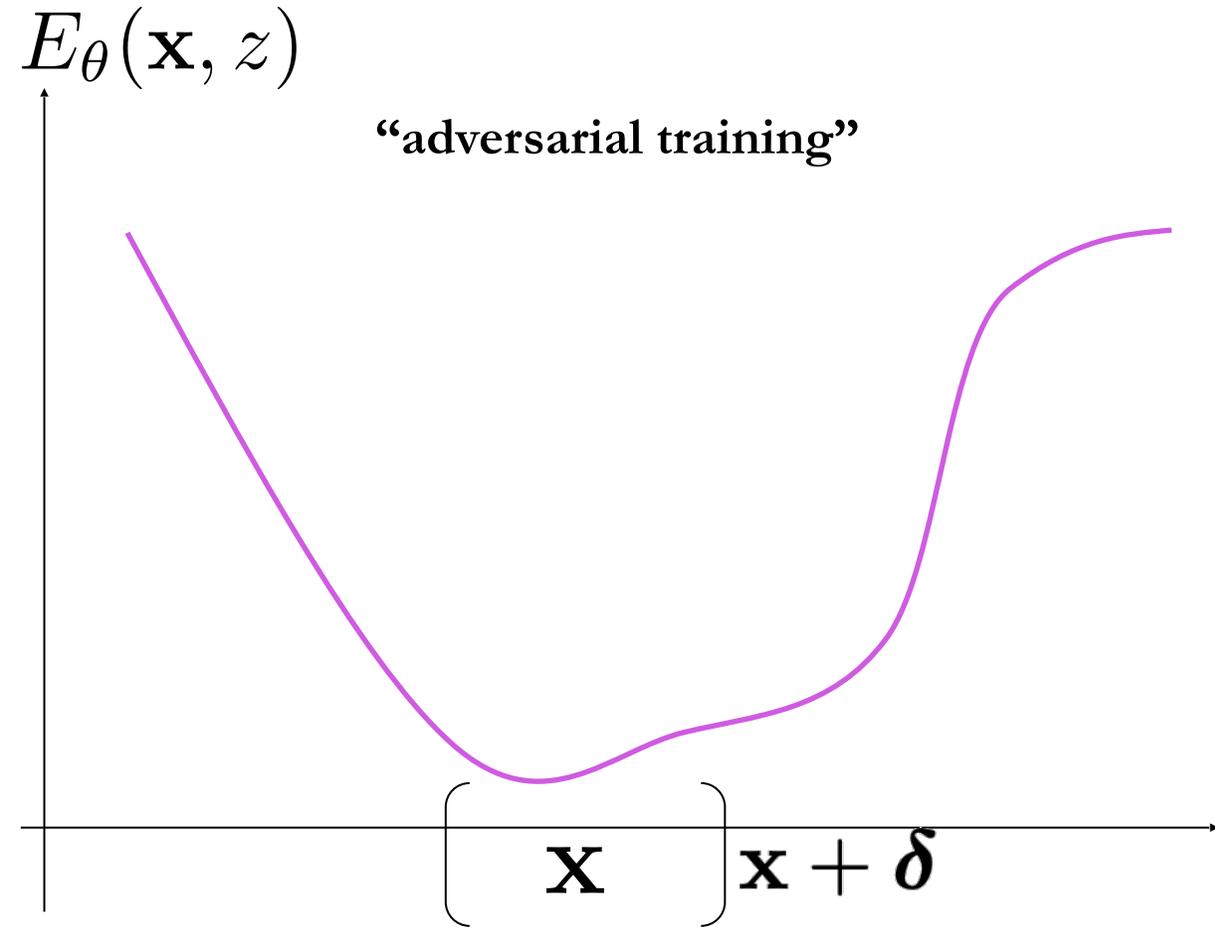


Why Robust Models behave as Generatives?



Why Robust Models behave as Generatives?





Energy $E_{\theta}(\mathbf{x})$ vs attack strength



Energy $E_{\theta}(\mathbf{x})$ vs attack strength



Finding 1: Untargeted attacks* decrease $E_{\theta}(\mathbf{x})$ thus increase $p_{\theta}(\mathbf{x})$

* Untargeted attack = PGD (Projected Gradient Descent) attack

Energy $E_{\theta}(\mathbf{x})$ vs attack strength



Finding 1: Untargeted attacks* decrease $E_{\theta}(\mathbf{x})$ thus increase $p_{\theta}(\mathbf{x})$

In other words, untargeted attacks finds points with:

- High energy $E_{\theta}(\mathbf{x}, z)$ thus low $p_{\theta}(\mathbf{x}, z)$

Fool the classifier (known)

* Untargeted attack = PGD (Projected Gradient Descent) attack

Energy $E_{\theta}(\mathbf{x})$ vs attack strength



Finding 1: Untargeted attacks* decrease $E_{\theta}(\mathbf{x})$ thus increase $p_{\theta}(\mathbf{x})$

In other words, untargeted attacks finds points with:

- High energy $E_{\theta}(\mathbf{x}, z)$ thus low $p_{\theta}(\mathbf{x}, z)$
- Low $E_{\theta}(\mathbf{x})$, highly likely for the model $p_{\theta}(\mathbf{x})$

Fool the classifier (known)

New adversarial points are more likely to exist than the natural data points! (less known)

* Untargeted attack = PGD (Projected Gradient Descent) attack

Energy $E_{\theta}(\mathbf{x})$ vs attack strength



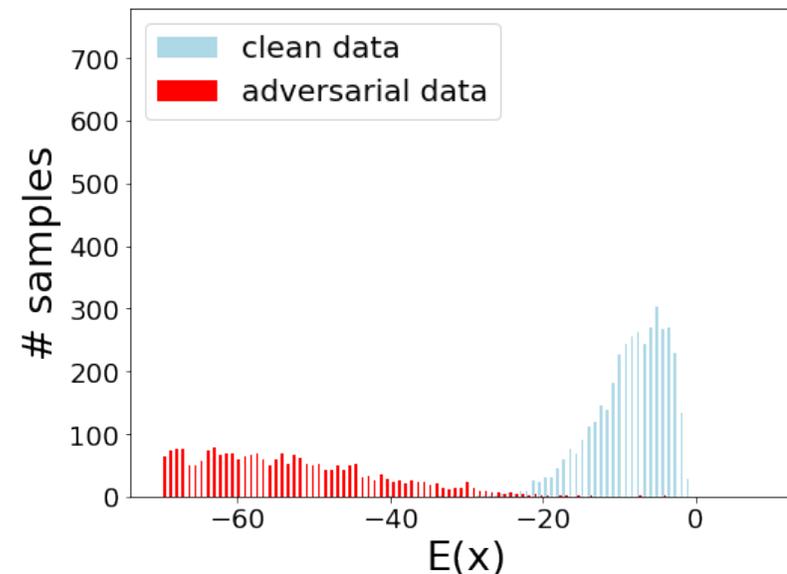
Finding 1: Untargeted attacks* decrease $E_{\theta}(\mathbf{x})$ thus increase $p_{\theta}(\mathbf{x})$

In other words, untargeted attacks finds points with:

- High energy $E_{\theta}(\mathbf{x}, z)$ thus low $p_{\theta}(\mathbf{x}, z)$
- Low $E_{\theta}(\mathbf{x})$, highly likely for the model $p_{\theta}(\mathbf{x})$

Fool the classifier (known)

New adversarial points are more likely to exist than the natural data points! (less known)



* Untargeted attack = PGD (Projected Gradient Descent) attack

Energy $E_{\theta}(\mathbf{x})$ vs attack strength



Energy $E_{\theta}(\mathbf{x})$ vs attack strength

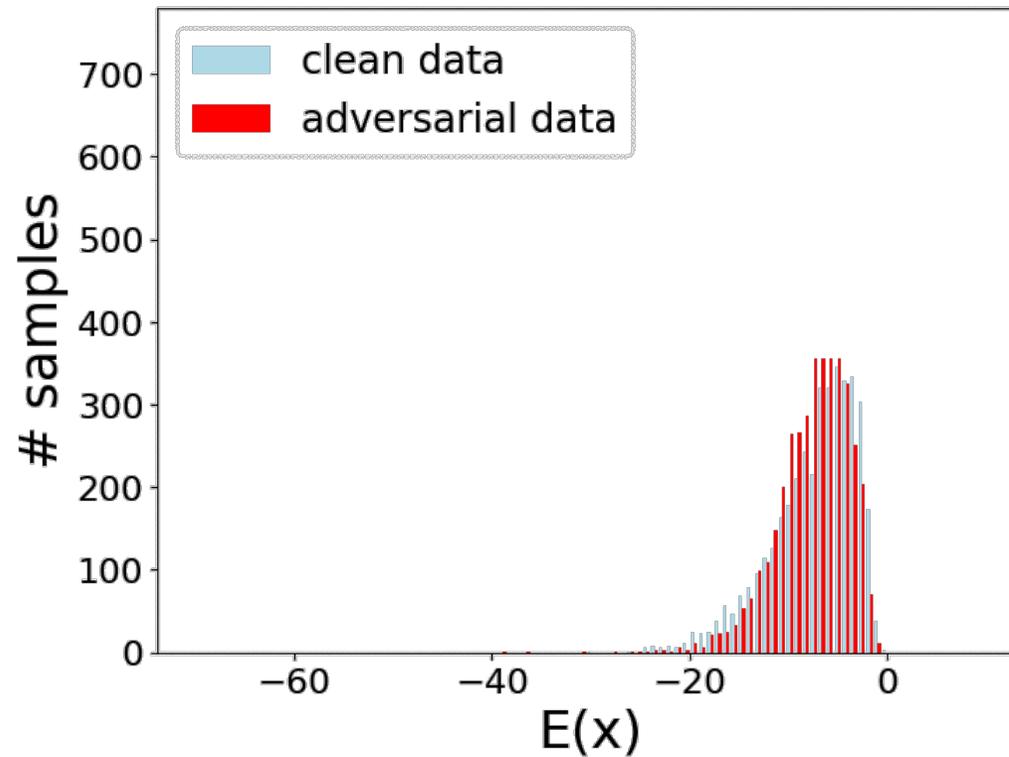


Finding 2: $E_{\theta}(\mathbf{x})$ decreases as the attack “strength” increases

Energy $E_\theta(\mathbf{x})$ vs attack strength



Finding 2: $E_\theta(\mathbf{x})$ decreases as the attack “strength” increases

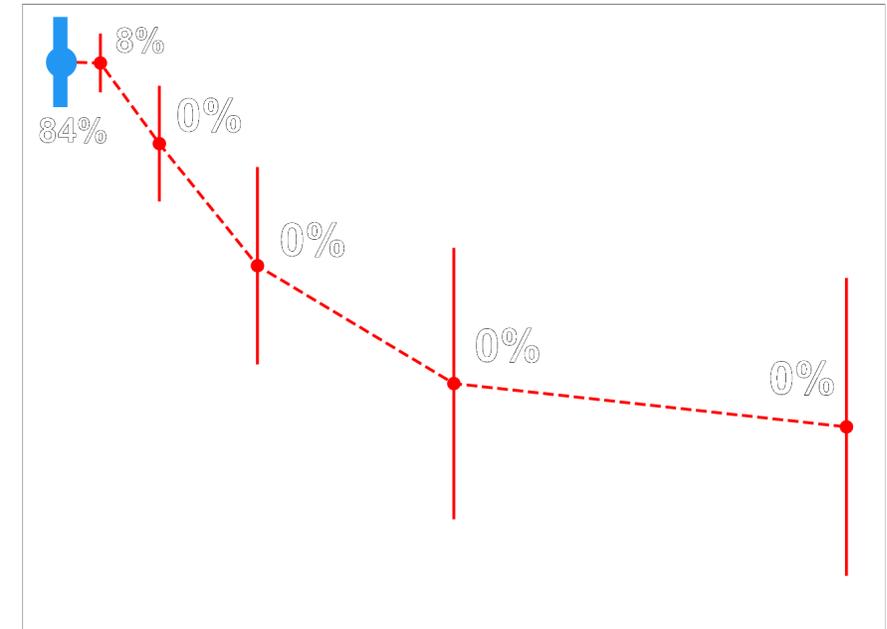
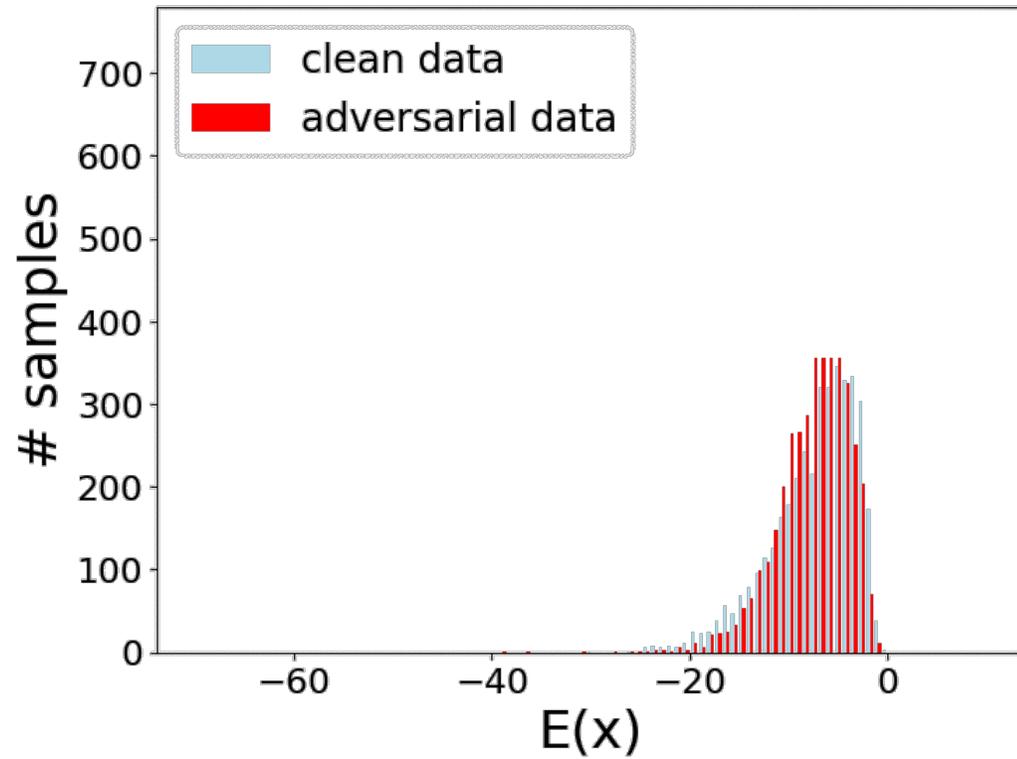


* attack strength= iterations in PGD

Energy $E_\theta(\mathbf{x})$ vs attack strength



Finding 2: $E_\theta(\mathbf{x})$ decreases as the attack “strength” increases

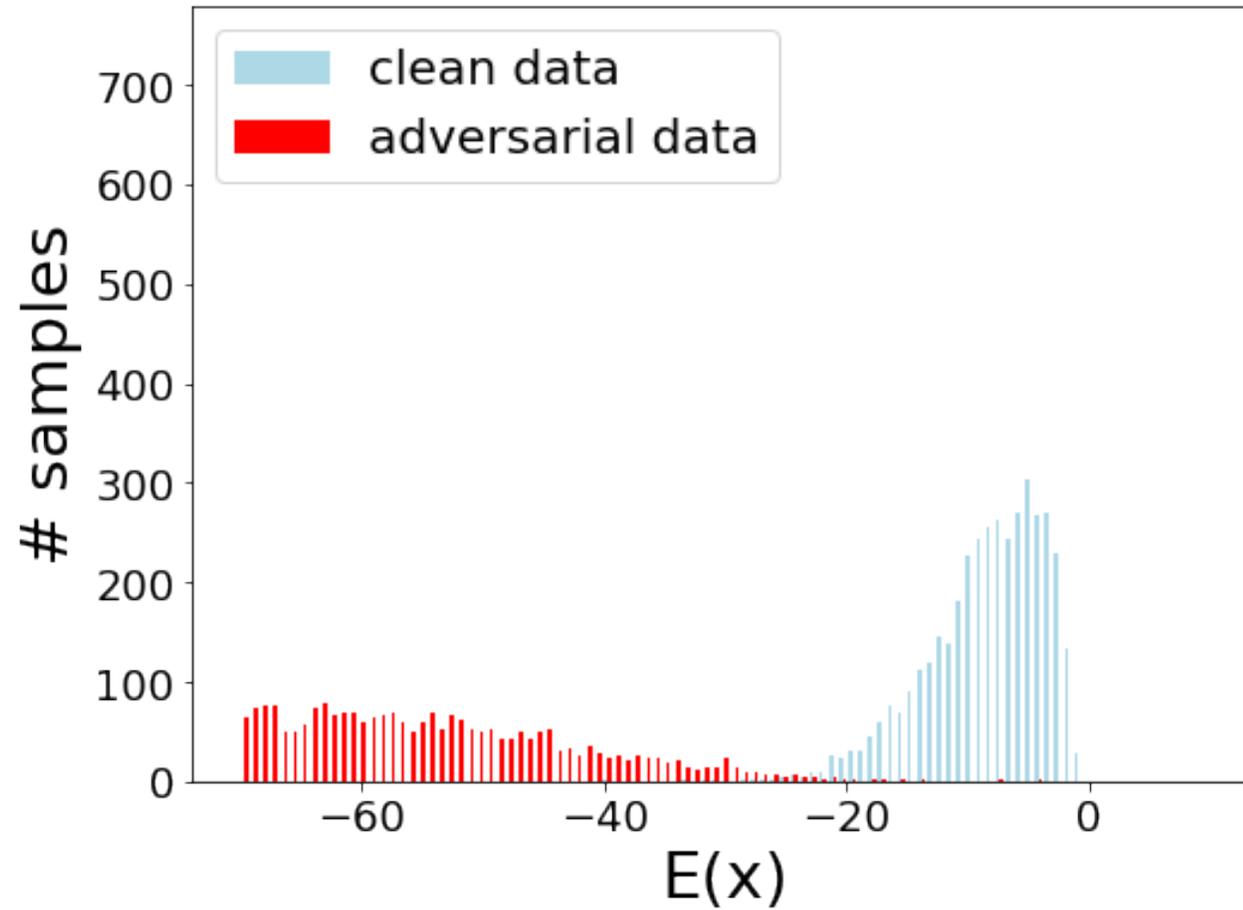


* attack strength= iterations in PGD

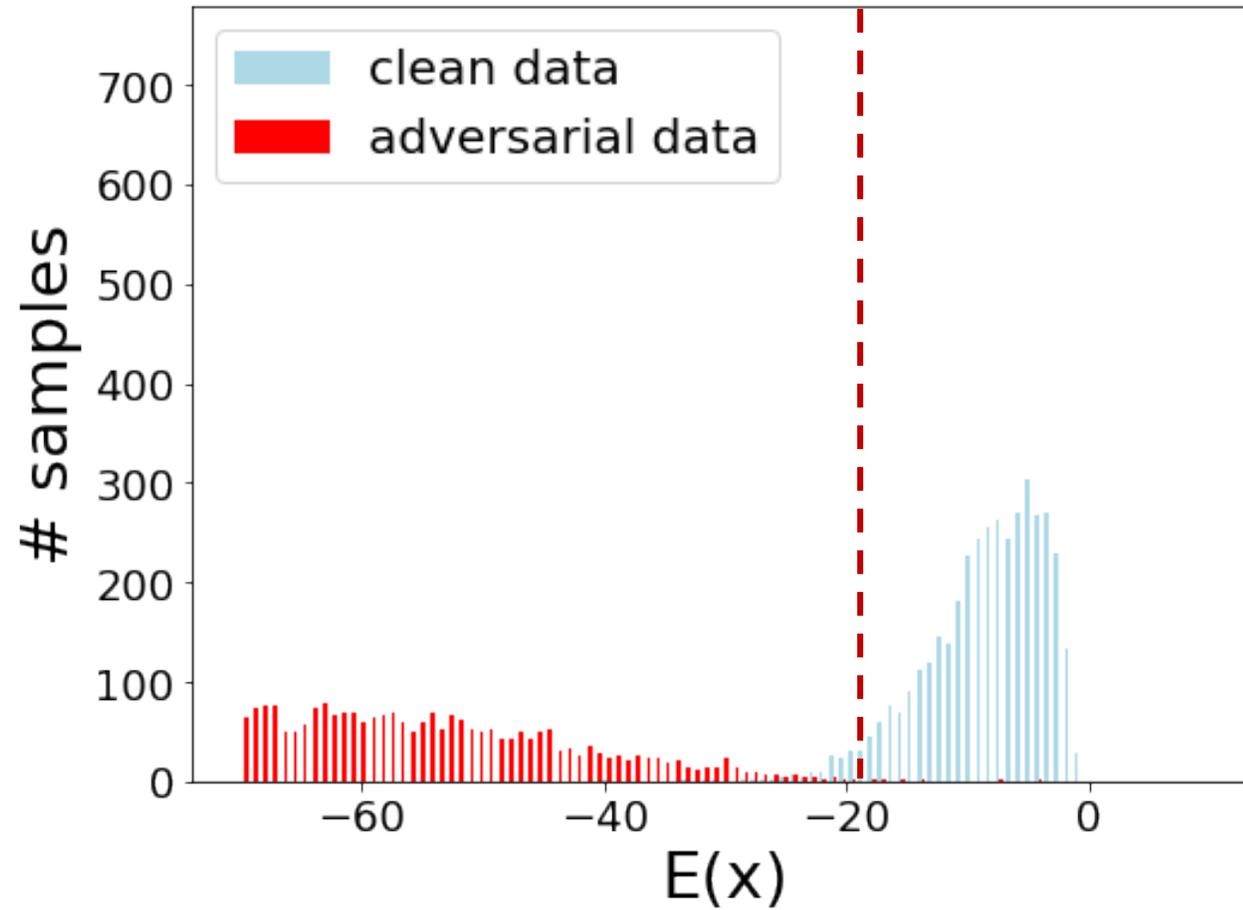
Efficient and Effective Attack Detector



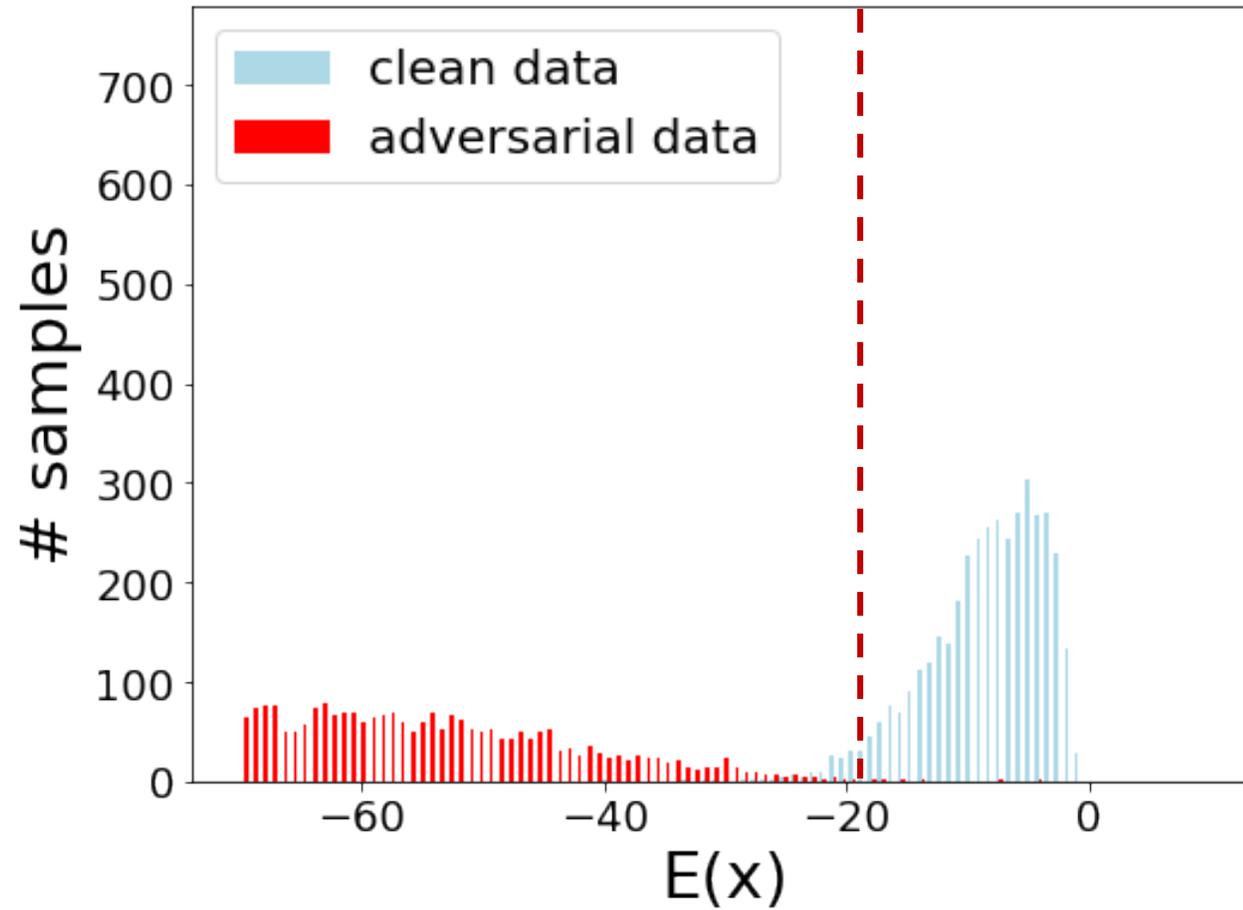
Efficient and Effective Attack Detector



Efficient and Effective Attack Detector



Efficient and Effective Attack Detector





Efficient and Effective Attack Detector



SAPIENZA
UNIVERSITÀ DI ROMA

Efficient and Effective Attack Detector



Dataset	Defense	Attack	DR	FPR
imagenette [8]	Energy (ResNet10)	PGD (8)	98.24	1.37
		PGD (16)	99.6	0.00

Efficient and Effective Attack Detector



Dataset	Defense	Attack	DR	FPR
imagenette [8]	Energy	PGD (8)	98.24	1.37
	(ResNet10)	PGD (16)	99.6	0.00

Dataset	Defense	Attack	DR	FPR
imagenette [8]	Energy (ResNet10)	PGD (8)	98.24	1.37
		PGD (16)	99.6	0.00
CIFAR-10 [17]	Energy (ResNet10)	PGD (8)	98.38	1.62
		APGD (8)	85.45	1.19
	KD+BU [9]	PGD (8)	92.27	0.96
	LID [21]	PGD (8)	94.39	1.81

Dataset	Defense	Attack	DR	FPR
imagenette [8]	Energy	PGD (8)	98.24	1.37
	(ResNet10)	PGD (16)	99.6	0.00
CIFAR-10 [17]	Energy	PGD (8)	98.38	1.62
	(ResNet10)	APGD (8)	85.45	1.19
	KD+BU [9]	PGD (8)	92.27	0.96
	LID [21]	PGD (8)	94.39	1.81

Dataset	Defense	Attack	DR	FPR
imagenette [8]	Energy	PGD (8)	98.24	1.37
	(ResNet10)	PGD (16)	99.6	0.00
CIFAR-10 [17]	Energy	PGD (8)	98.38	1.62
	(ResNet10)	APGD (8)	85.45	1.19
	KD+BU [9]	PGD (8)	92.27	0.96
	LID [21]	PGD (8)	94.39	1.81

Detector may suffer from: (1) targeted attacks (2) AutoAttack

Efficient and Effective Attack Detector



Can we bypass the detector?



High-Energy PGD



SAPIENZA
UNIVERSITÀ DI ROMA

High-Energy PGD



Fool the classifier yet keep energy like natural data

High-Energy PGD



Fool the classifier yet keep energy like natural data

$$\arg \max_{\delta} \left[\mathcal{L}(\theta(\mathbf{x} + \delta), y) + \lambda E_{\theta}(\mathbf{x} + \delta) \right]$$

High-Energy PGD



Fool the classifier yet keep energy like natural data

High-Energy PGD



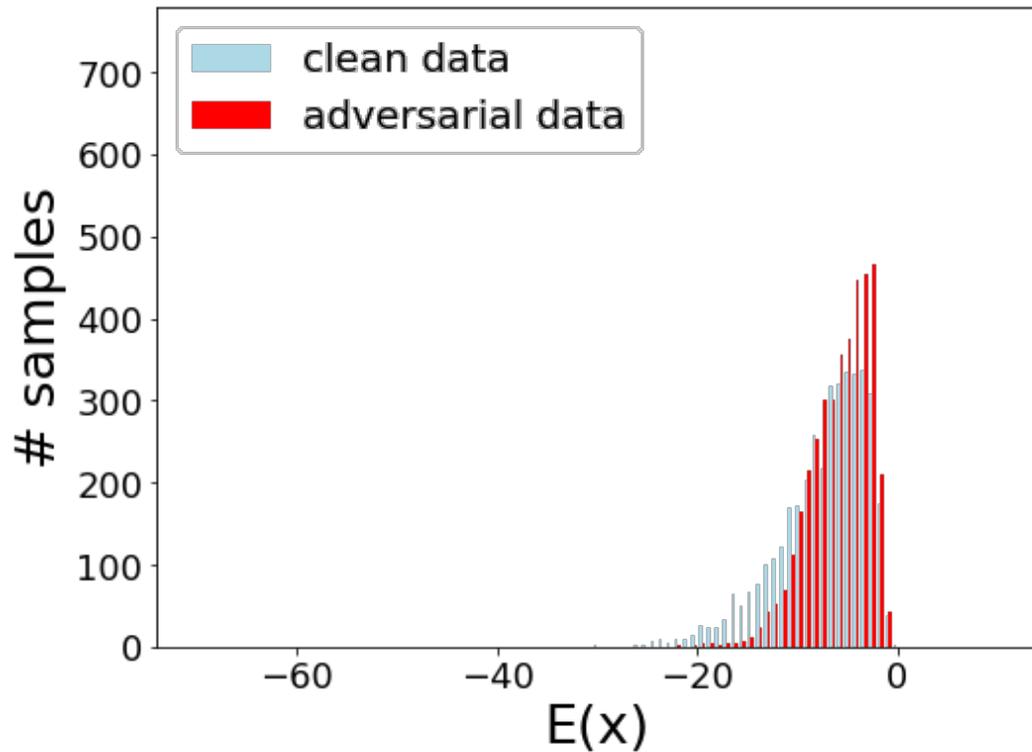
Fool the classifier yet keep energy like natural data

$$\mathbf{x}^* = \text{clip}_\epsilon \left[\mathbf{x}^* + \alpha \text{sign} \left[\nabla_{\mathbf{x}^*} \mathcal{L}(\boldsymbol{\theta}(\mathbf{x}^*), y) + \lambda E_{\boldsymbol{\theta}}(\mathbf{x}^*) \right] \right]$$

High-Energy PGD

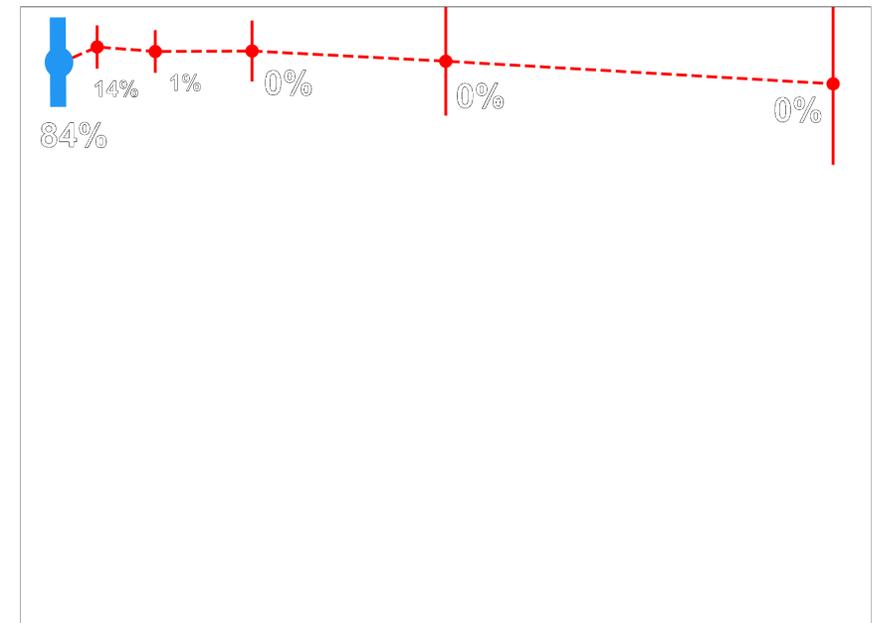
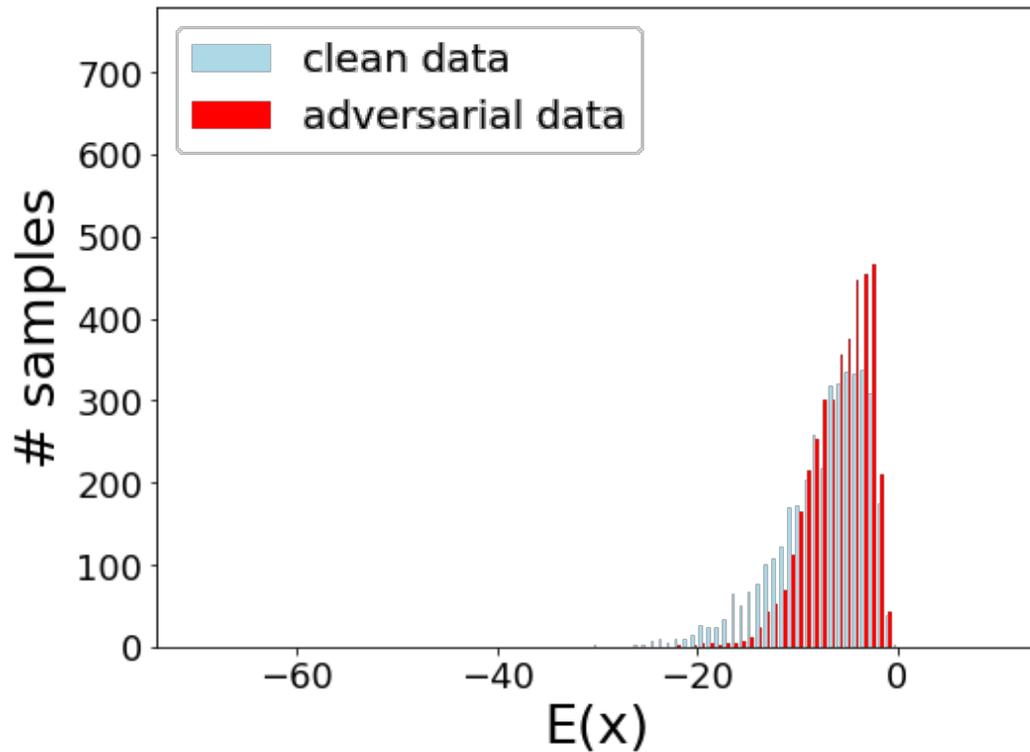
Fool the classifier yet keep energy like natural data

$$\mathbf{x}^* = \text{clip}_\epsilon \left[\mathbf{x}^* + \alpha \text{sign} \left[\nabla_{\mathbf{x}^*} \mathcal{L}(\boldsymbol{\theta}(\mathbf{x}^*), y) + \lambda E_{\boldsymbol{\theta}}(\mathbf{x}^*) \right] \right]$$



Fool the classifier yet keep energy like natural data

$$\mathbf{x}^* = \text{clip}_\epsilon \left[\mathbf{x}^* + \alpha \text{sign} \left[\nabla_{\mathbf{x}^*} \mathcal{L}(\boldsymbol{\theta}(\mathbf{x}^*), y) + \lambda E_{\boldsymbol{\theta}}(\mathbf{x}^*) \right] \right]$$



Exploring the Connection between Robust and Generative Models

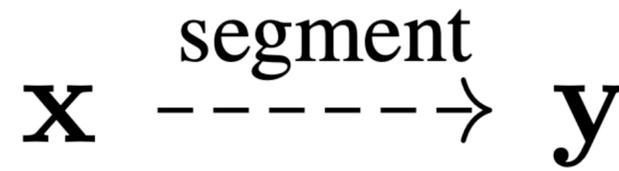
MAsk-Guided Image Synthesis by Inverting a Quasi-Robust Classifier [AAAI23]

Joint work: Mozhddeh Rouhsedaghat (USC)

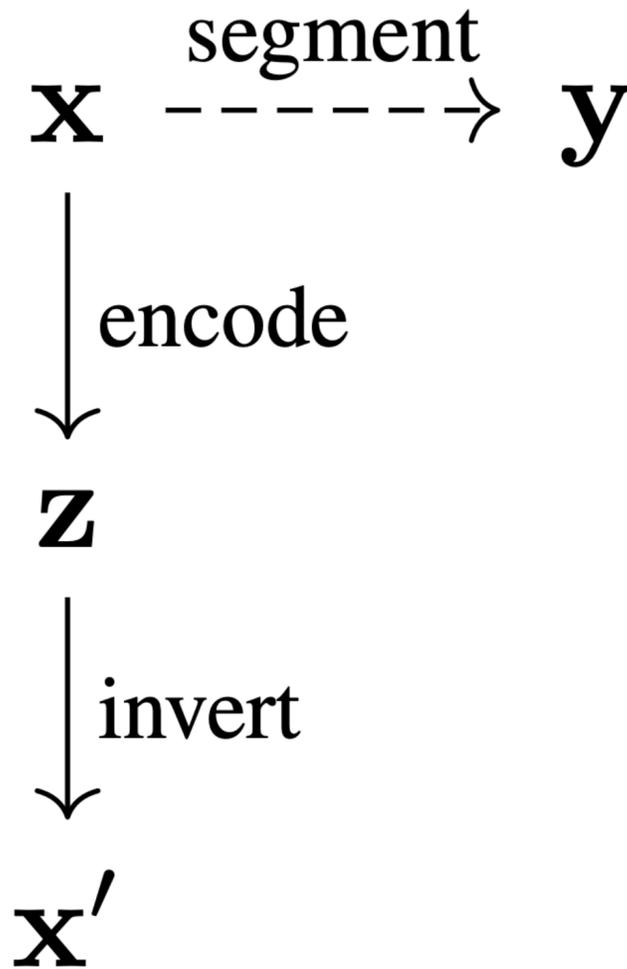
Masoud Monajatipoor (UCLA)

C.-C. Jay Kuo (USC)

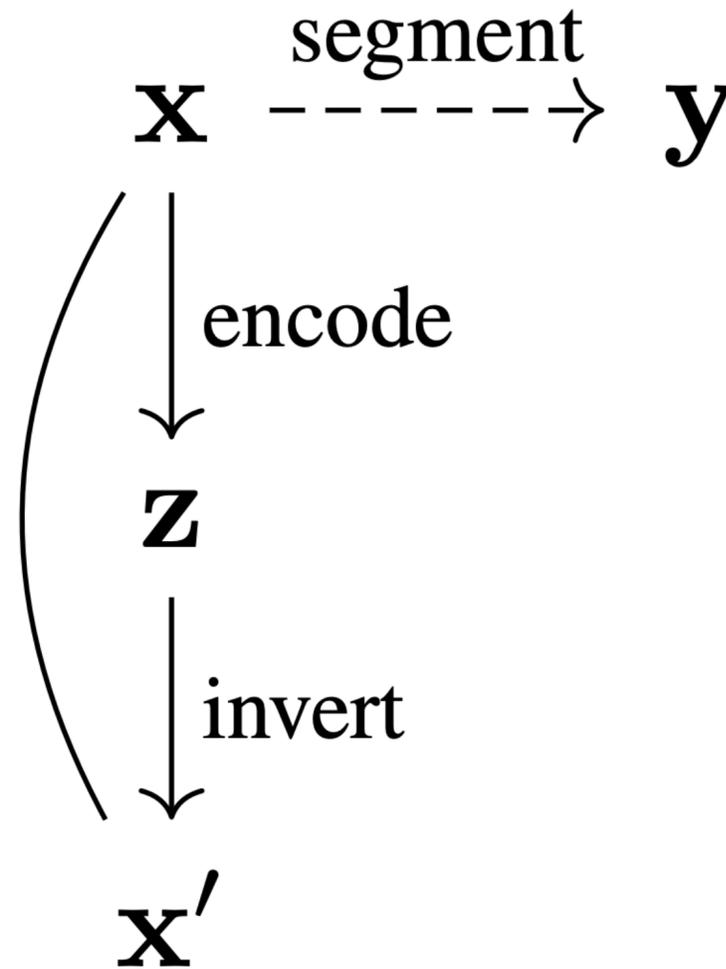
Method



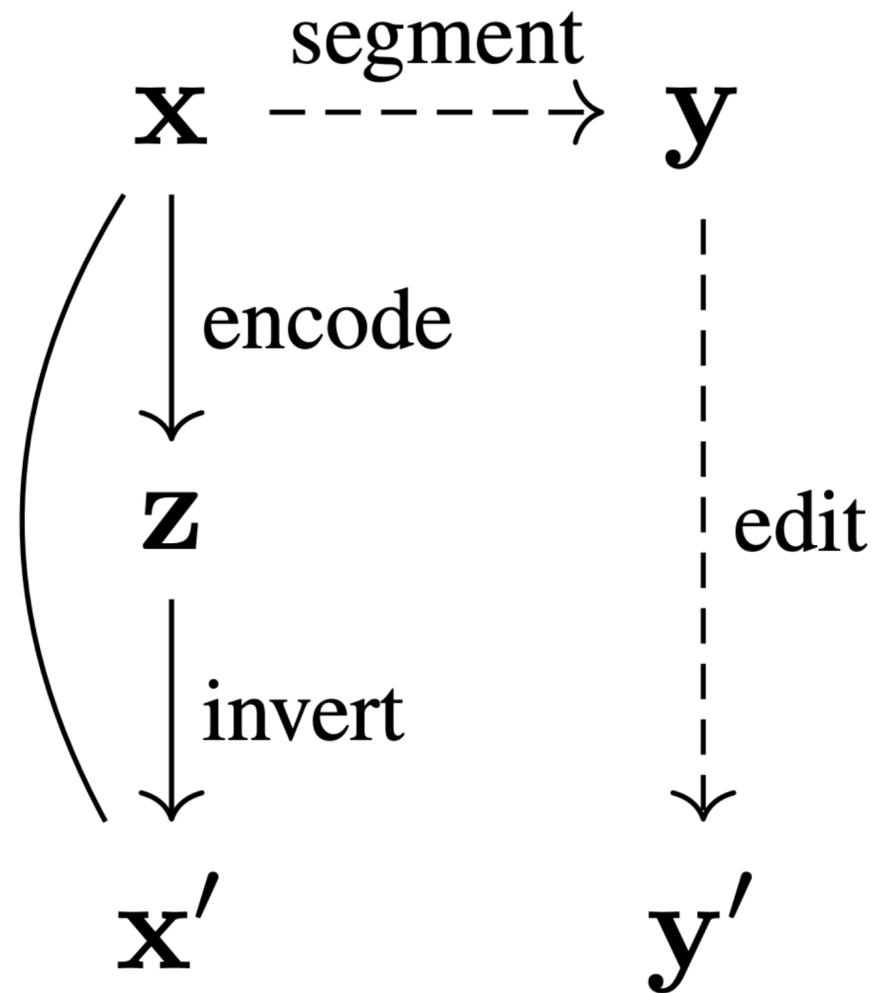
Method



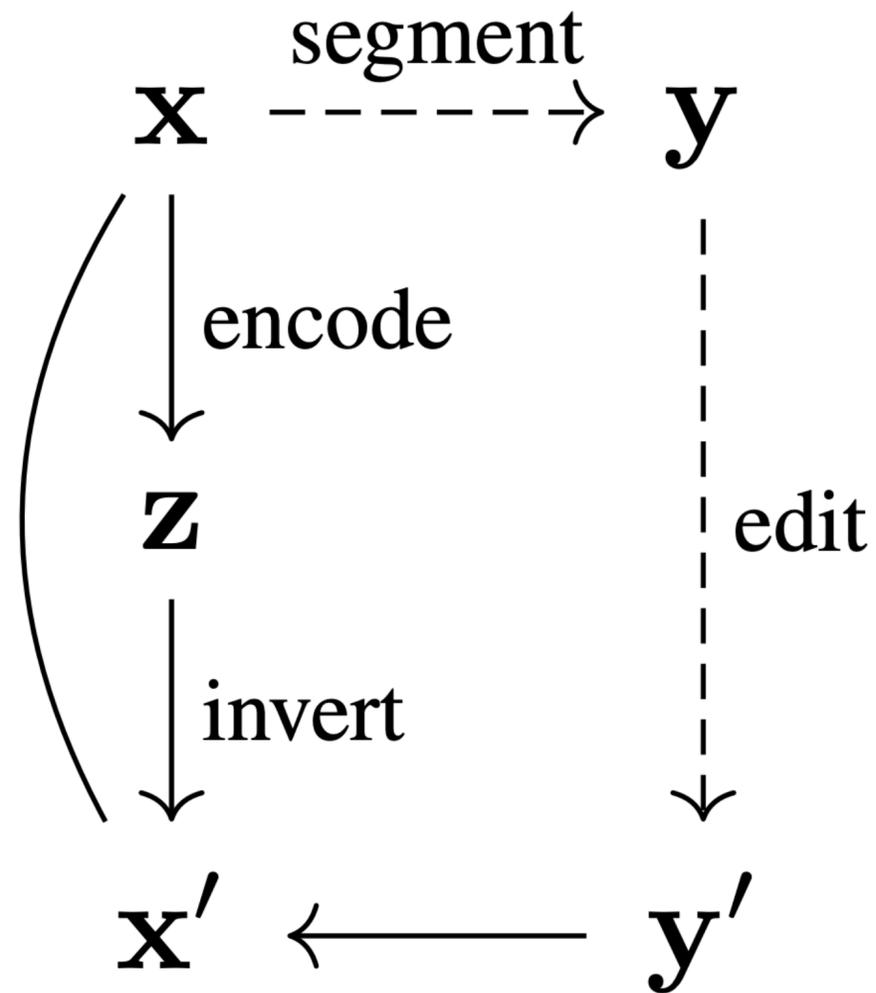
Method



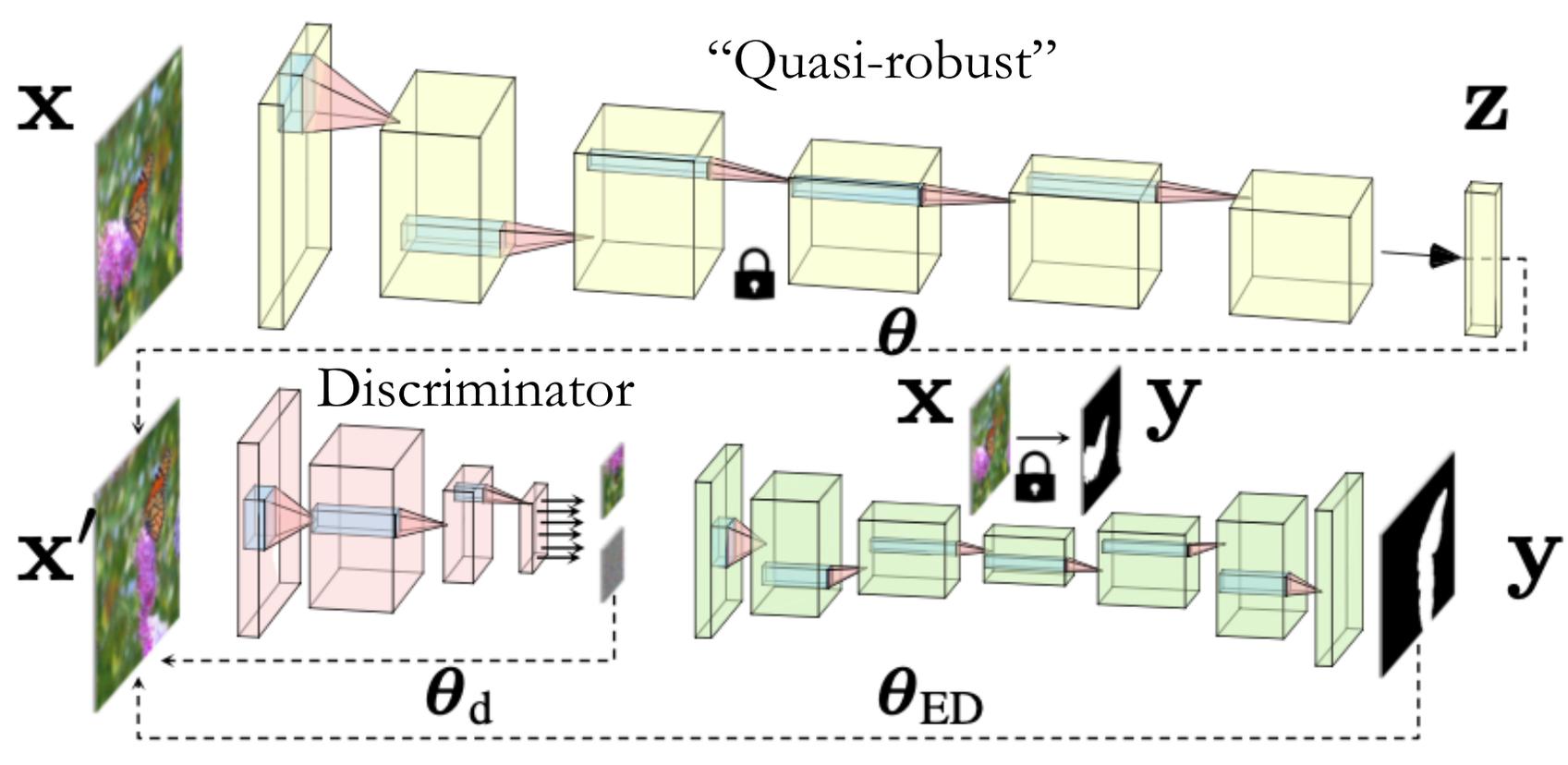
Method



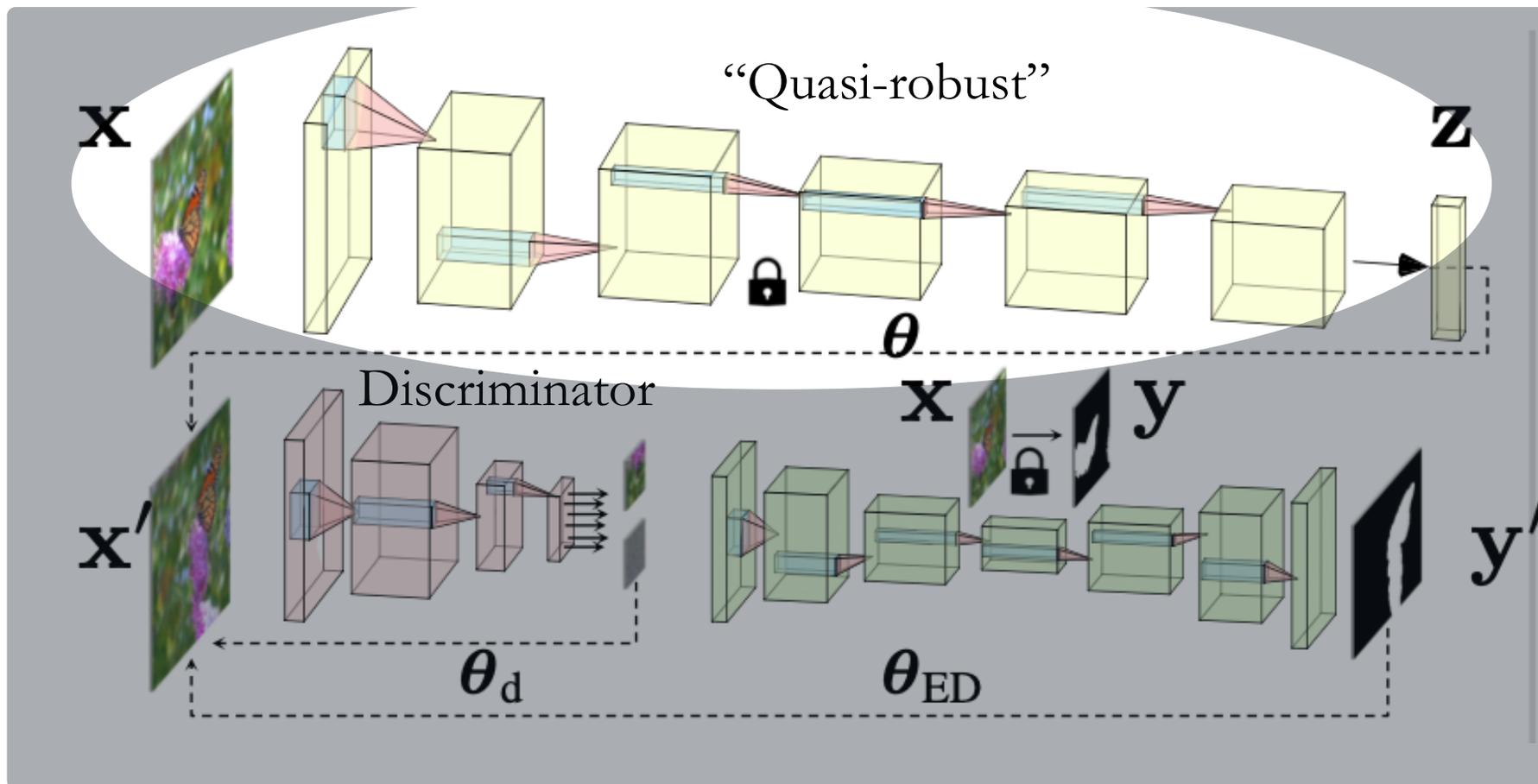
Method



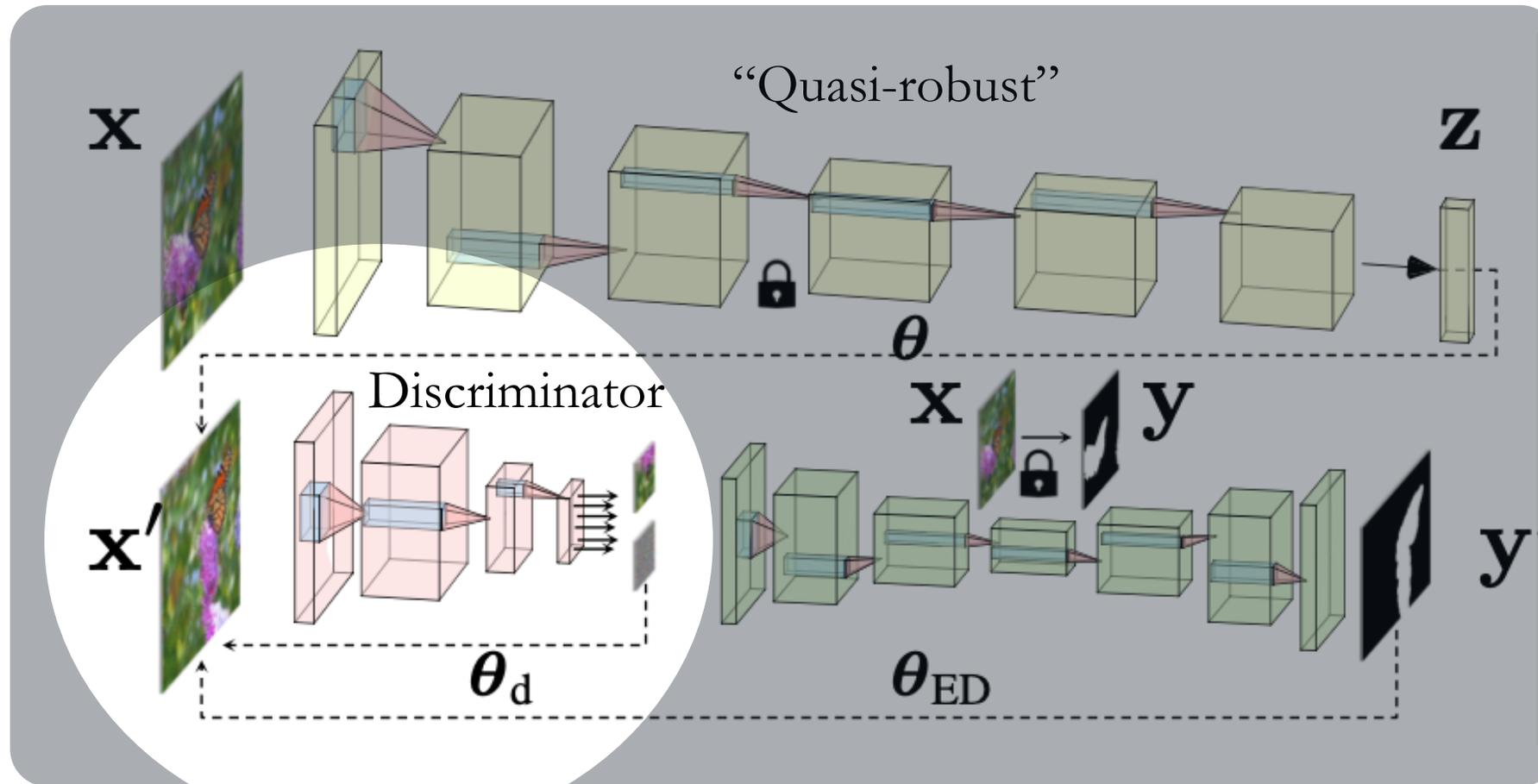
Manipulation Control



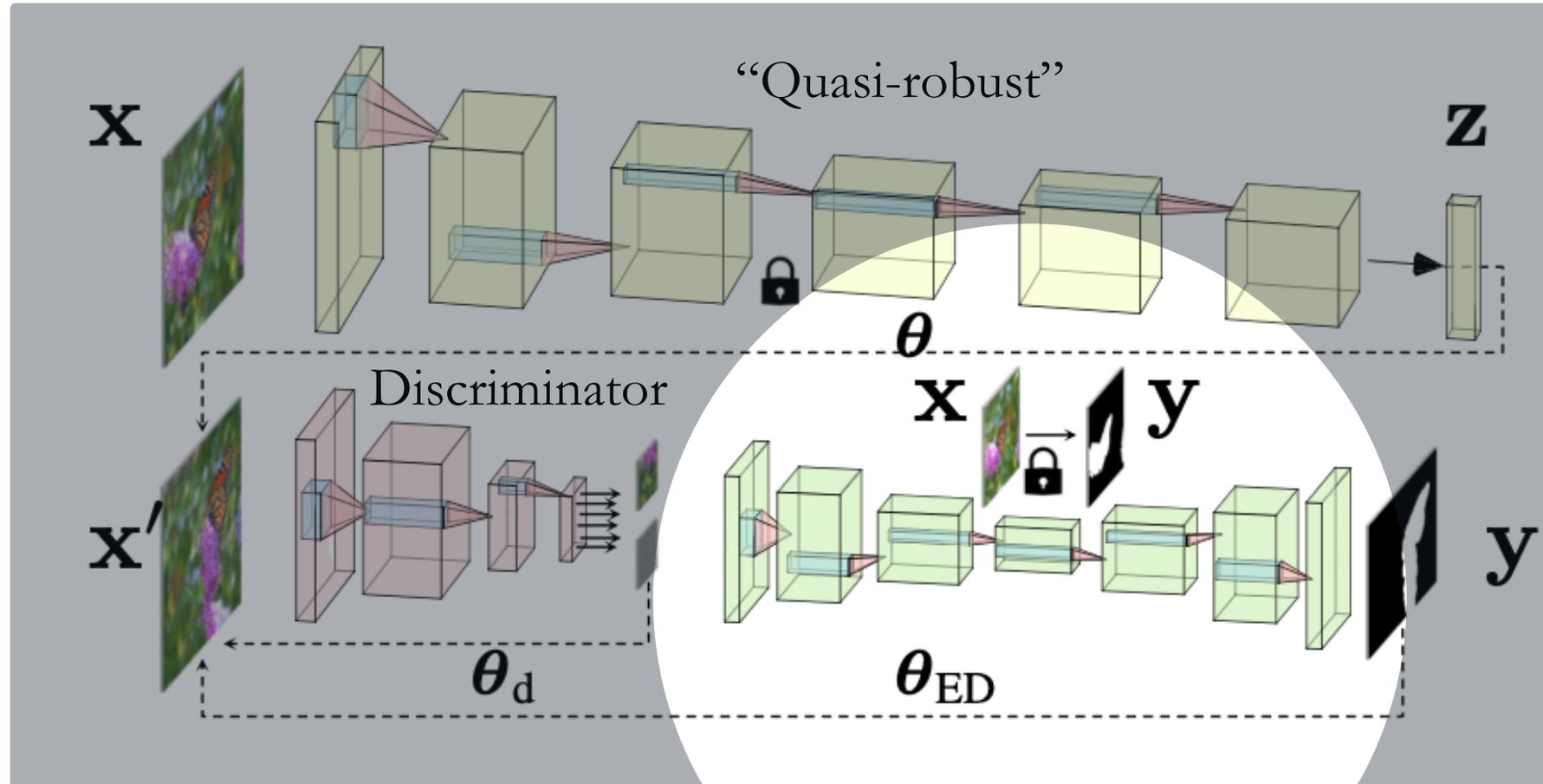
Manipulation Control



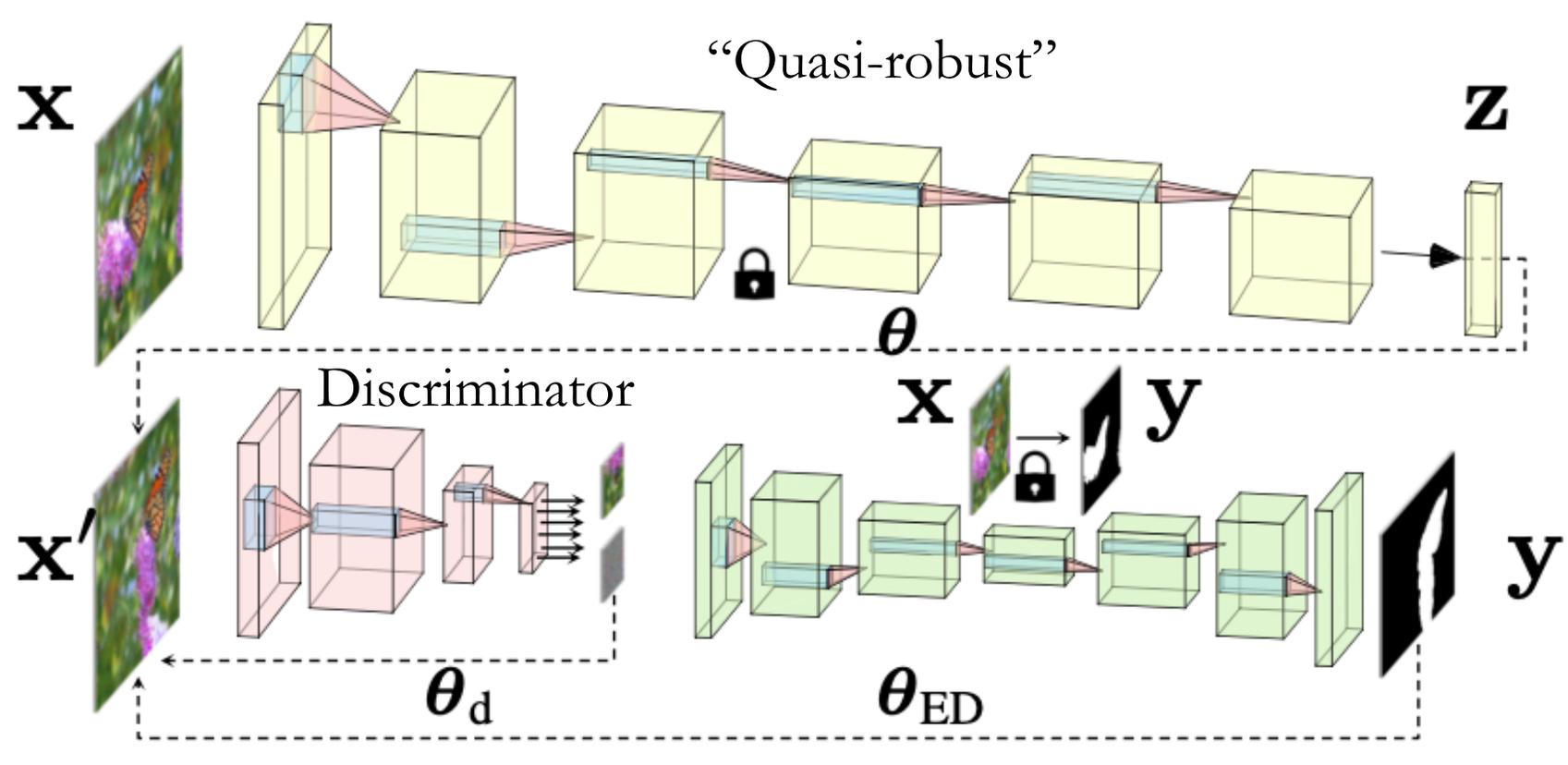
Manipulation Control



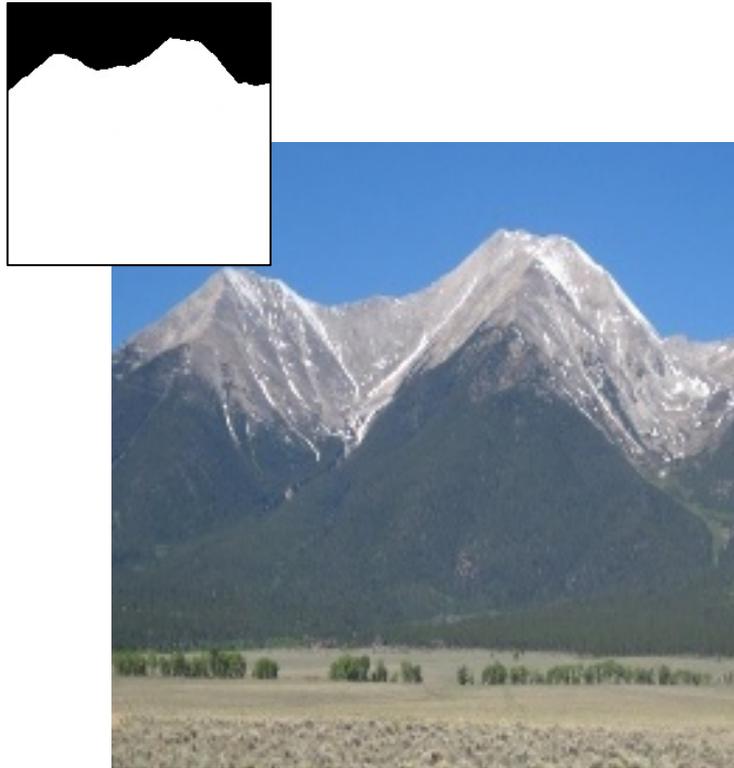
Manipulation Control



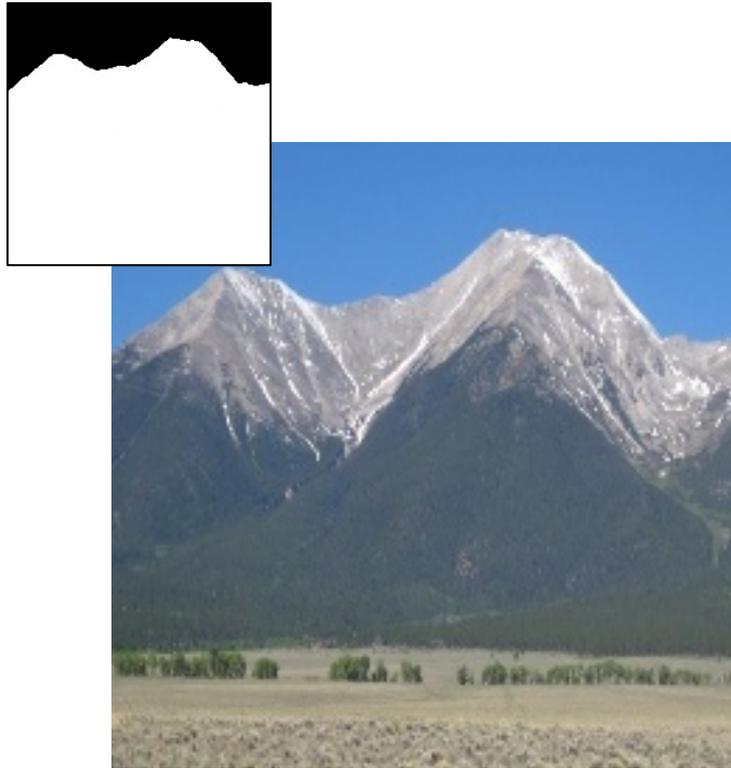
Manipulation Control



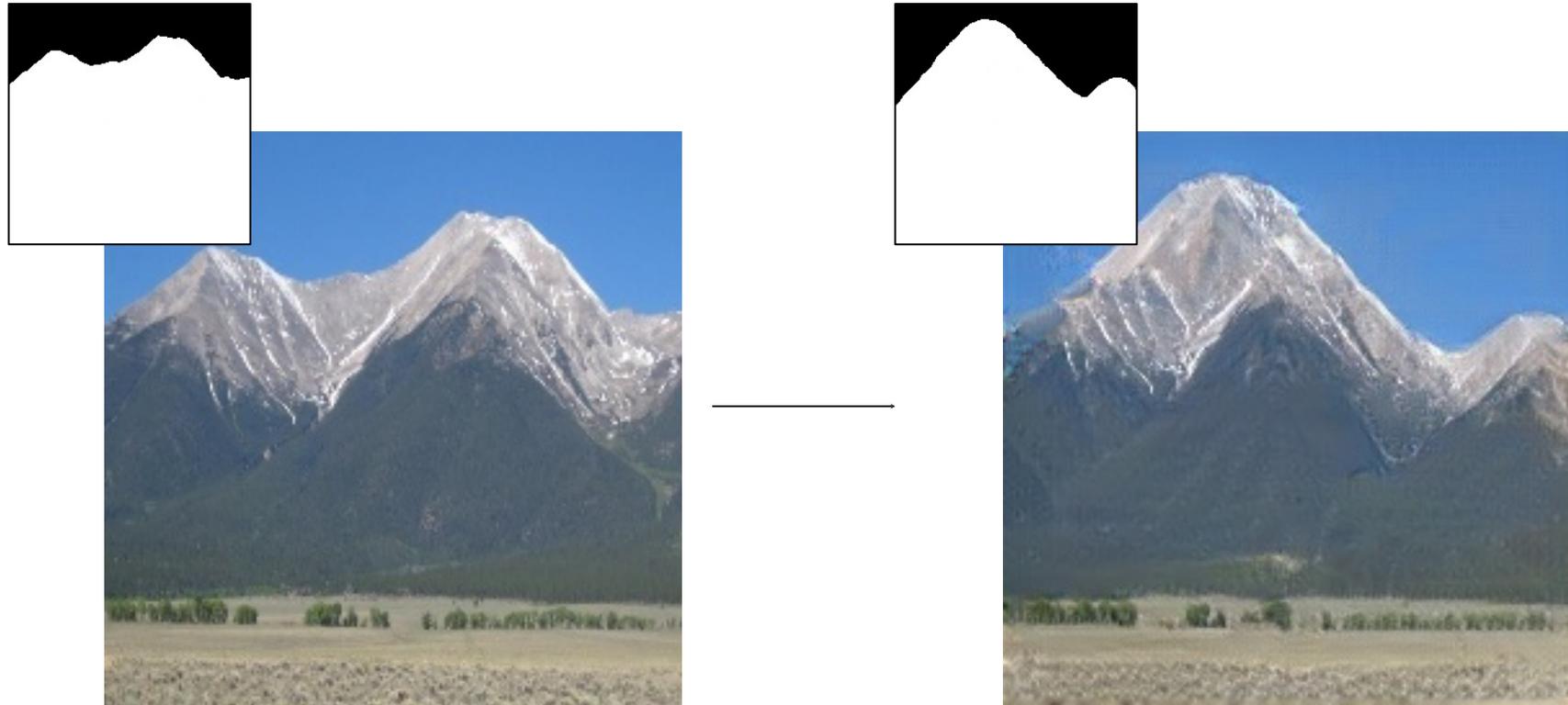
Manipulation Control – Non-Rigid Scene Deformation



Manipulation Control – Non-Rigid Scene Deformation



Manipulation Control – Non-Rigid Scene Deformation



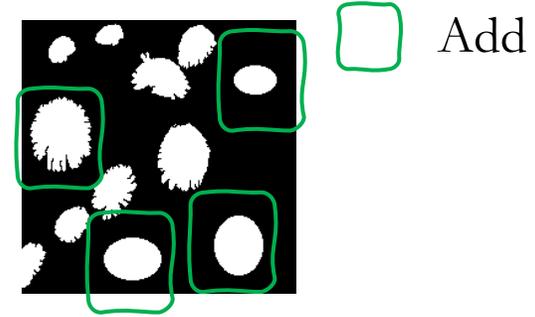
Manipulation Control – Copy/Move



Manipulation Control – Copy/Move



Manipulation Control – Copy/Move



Manipulation Control – Copy/Move

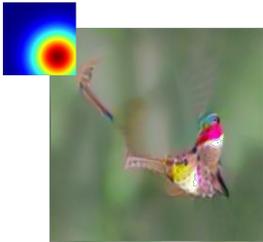


Qualitative Comparison

Input



IMAGINE
[4]



DEEPSIM
[5]



MAGIC
(Ours)



a)

Qualitative Comparison

Input



IMAGINE
[4]



DEEPSIM
[5]



MAGIC
(Ours)



a)

b)

Qualitative Comparison

Input



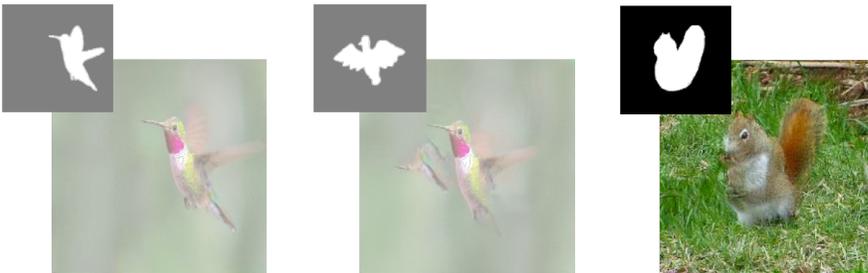
IMAGINE
[4]



DEEPSIM
[5]



MAGIC
(Ours)



a)

b)

c)

Qualitative Comparison

Input



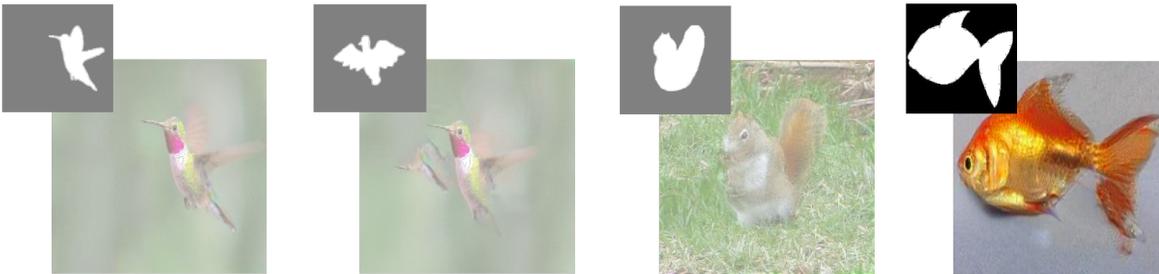
IMAGINE
[4]



DEEPSIM
[5]



MAGIC
(Ours)



a)

b)

c)

d)

Qualitative Comparison

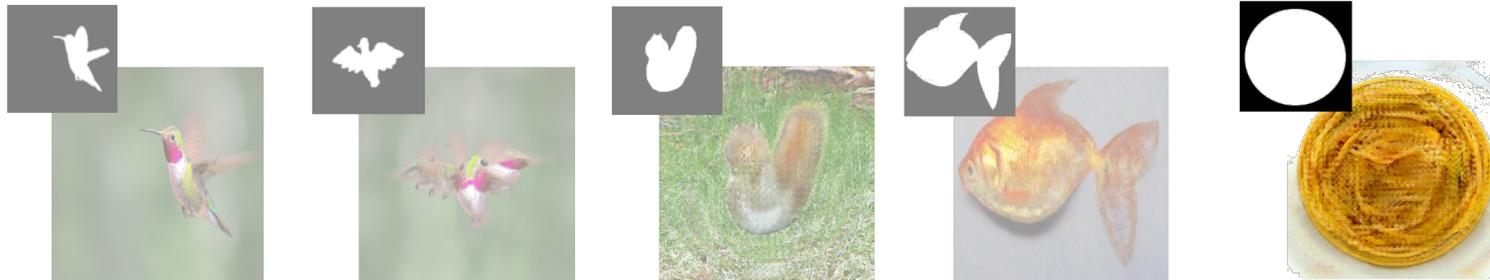
Input



IMAGINE
[4]



DEEPSIM
[5]



MAGIC
(Ours)



a)

b)

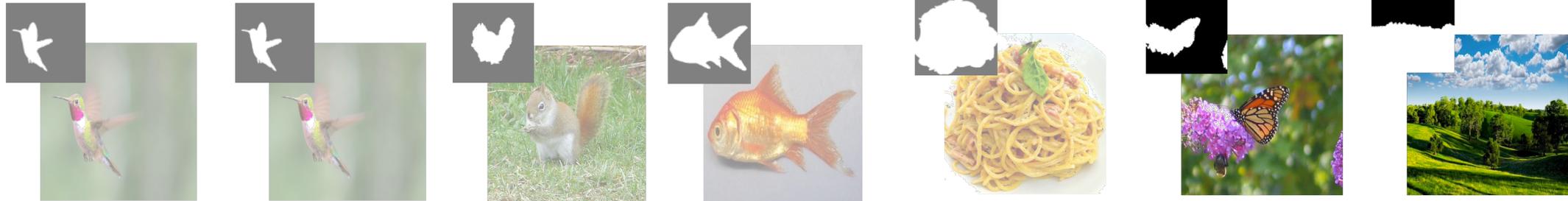
c)

d)

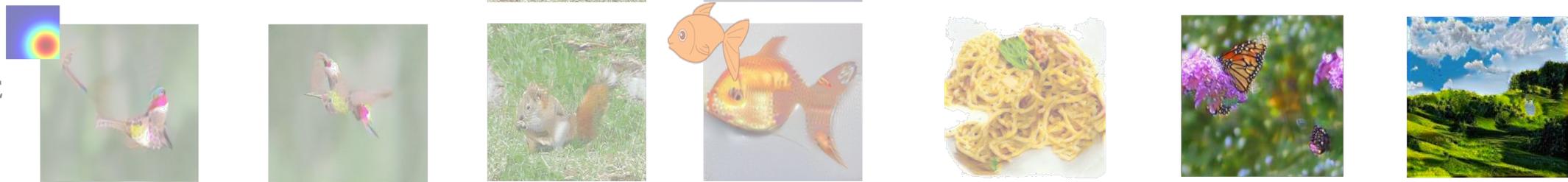
e)

Qualitative Comparison

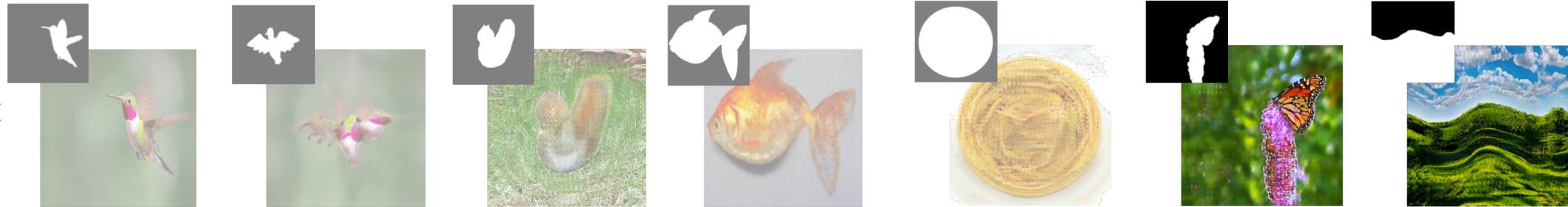
Input



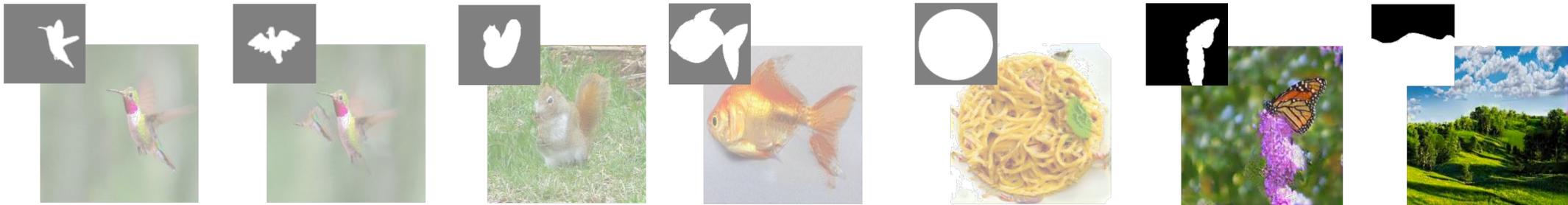
IMAGINE
[4]



DEEPSIM
[5]



MAGIC
(Ours)



a)

b)

c)

d)

e)

f)

g)



Future Work



SAPIENZA
UNIVERSITÀ DI ROMA

Future Work



**Better analyze
High-Energy PGD**

**Investigate same but
for **targeted** attacks**

Future Work



**Better analyze
High-Energy PGD**

Investigate same but
for **targeted** attacks

**Investigate Hybrid
Generative-Discriminative
Models**

Future Work



Thank you!