

Go to www.menti.com and use the code **53 39 35**

CS230: Lecture 5

Attacking Networks with Adversarial Examples

-

Generative Adversarial Networks

Kian Katanforoosh

Today's outline

- I. Attacking NNs with Adversarial Examples
- II. Generative Adversarial Networks

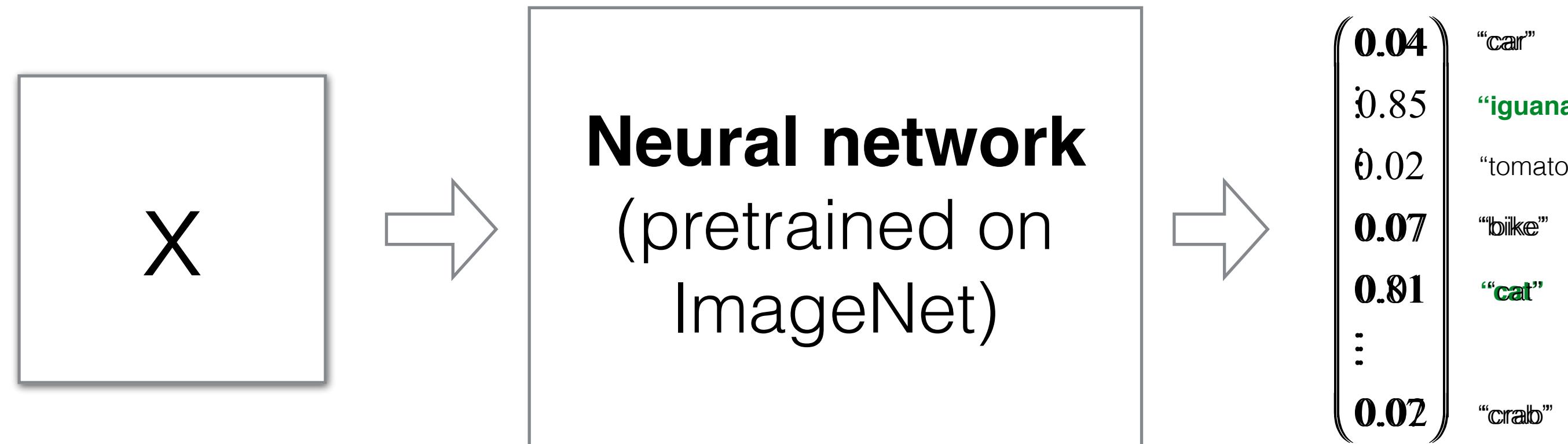
I. Adversarial examples

Discovery (2014): several machine learning models, including state-of-the-art neural networks, are vulnerable to adversarial examples

- A. How to build adversarial examples and attack a network?
- B. Examples
- C. How to defend against adversarial examples?

I. A. How to build adversarial examples and attack a network?

Goal: Given a pretrained network on ImageNet, find an example that is not a iguana but will be classify as an iguana.



1. Rephrasing what we want:

Find x such that: $\hat{y}(x) = y_{iguana} = \begin{pmatrix} 0 \\ 1 \\ \vdots \\ 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$

2. Defining the loss function

$$L(\hat{y}, y) = \frac{1}{2} \left\| \hat{y}(W, b, x) - y_{iguana} \right\|_2^2$$

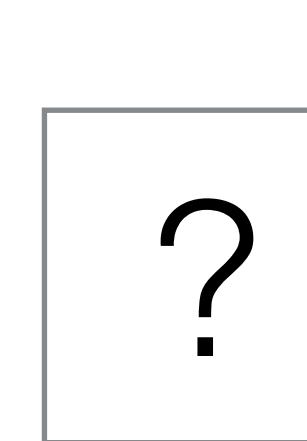
3. Optimize the image

After many iterations

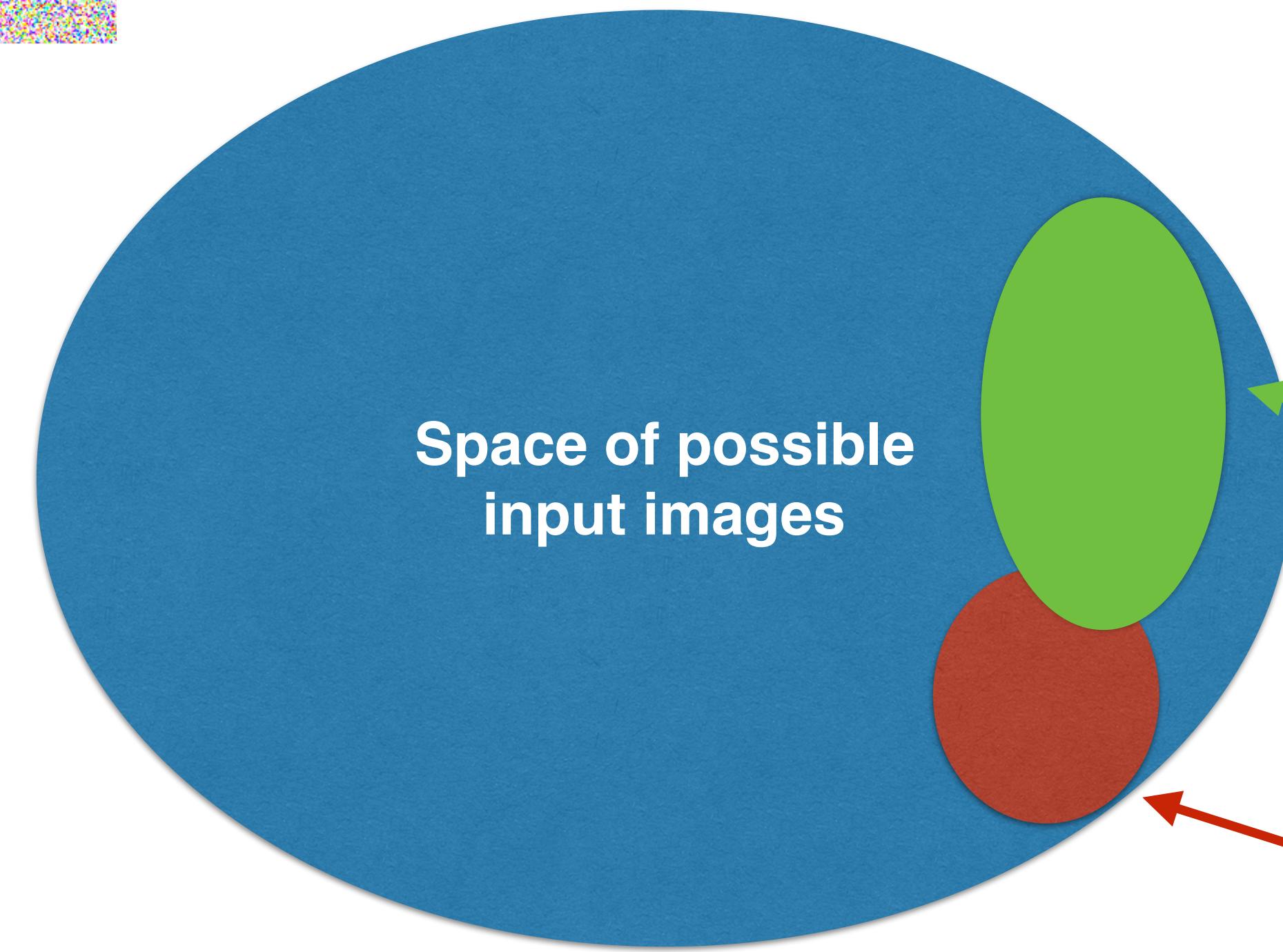
$$\frac{\partial L}{\partial x} \quad x = x - \alpha \frac{\partial L}{\partial x}$$

I. A. How to build adversarial examples and attack a network?

Question: Will the learned image \mathbf{x} look like an **iguana**?



$$256^{32 \times 32 \times 3} \approx 10^{7400}$$

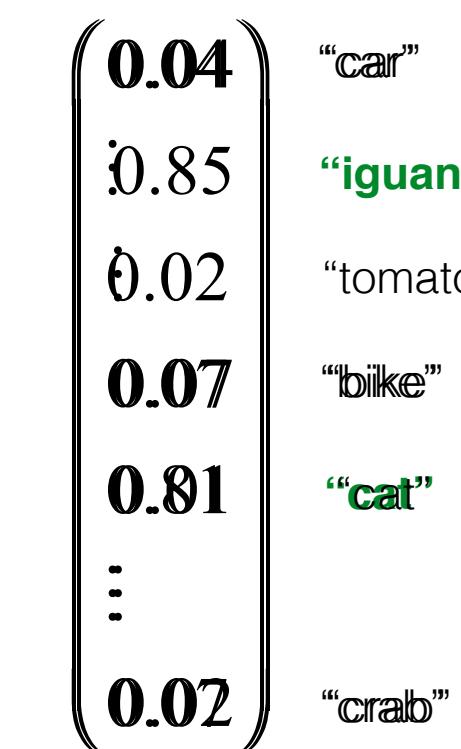
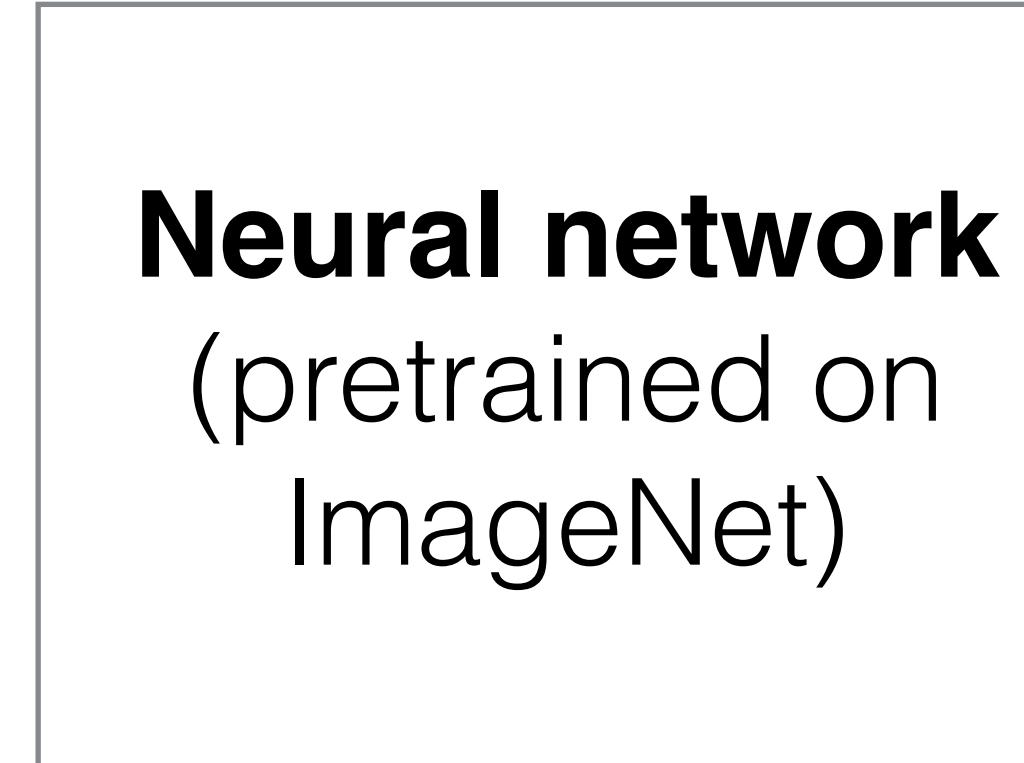
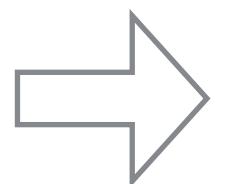
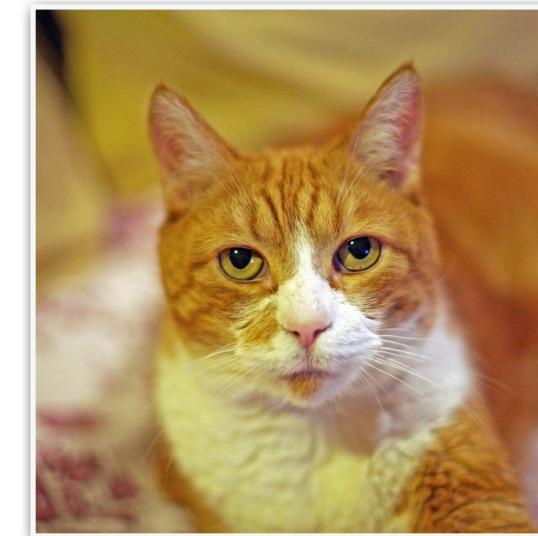


Space of images classified
as iguanas

Space of real images

I. A. How to build adversarial examples and attack a network?

Goal: Given a pretrained network on ImageNet, find an example that is a cat but will be classified as an iguana.



1. Rephrasing what we want:

Find x such that: $\hat{y}(x) = y_{iguana} = \begin{pmatrix} 0 \\ 1 \\ \vdots \\ 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$

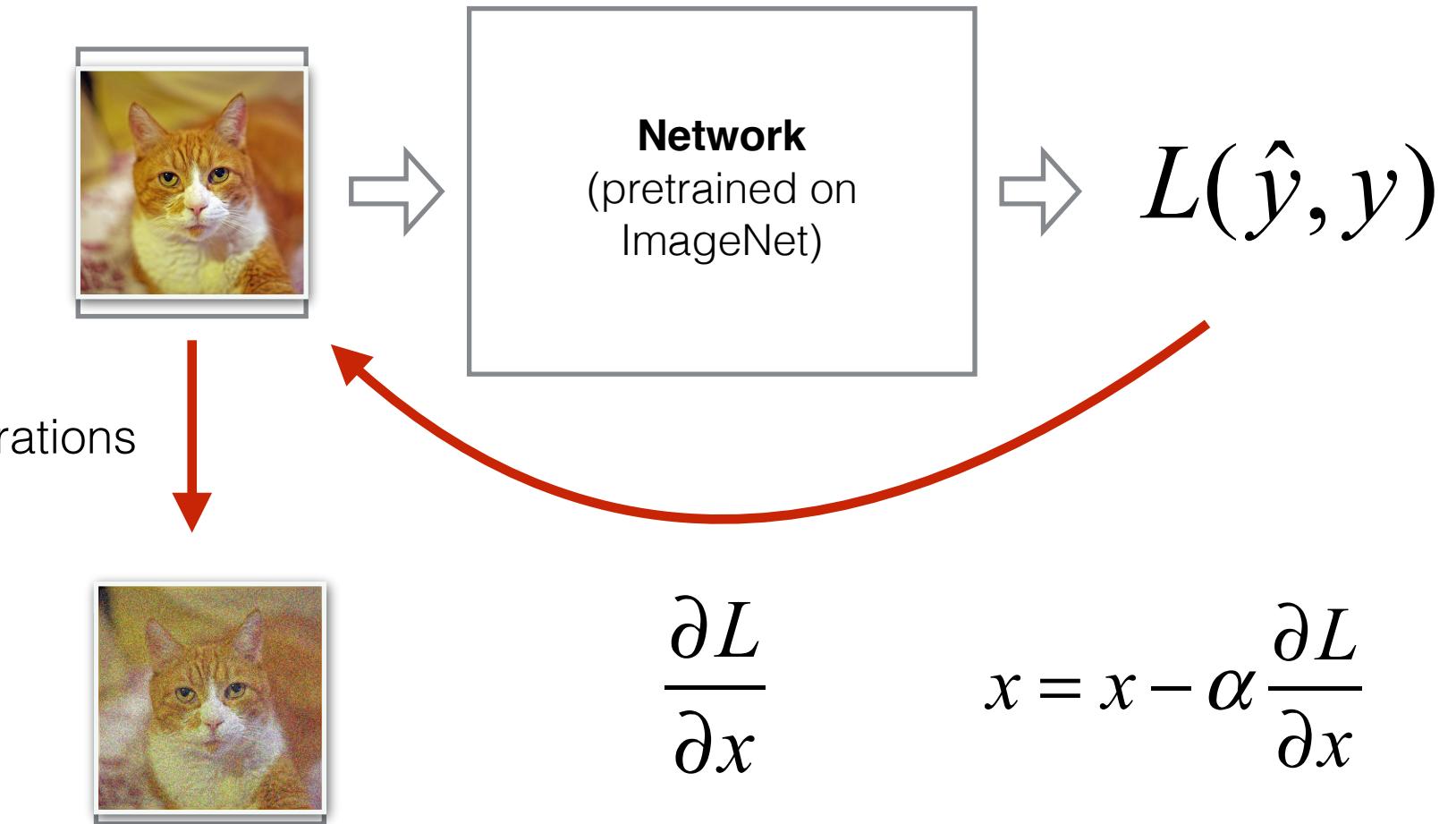
And: $x = x_{cat}$

2. Defining the loss function

$$L(\hat{y}, y) = \frac{1}{2} \left\| \hat{y}(W, b, x) - y_{iguana} \right\|_2^2 + \lambda \left\| x - x_{cat} \right\|_2^2$$

After many iterations

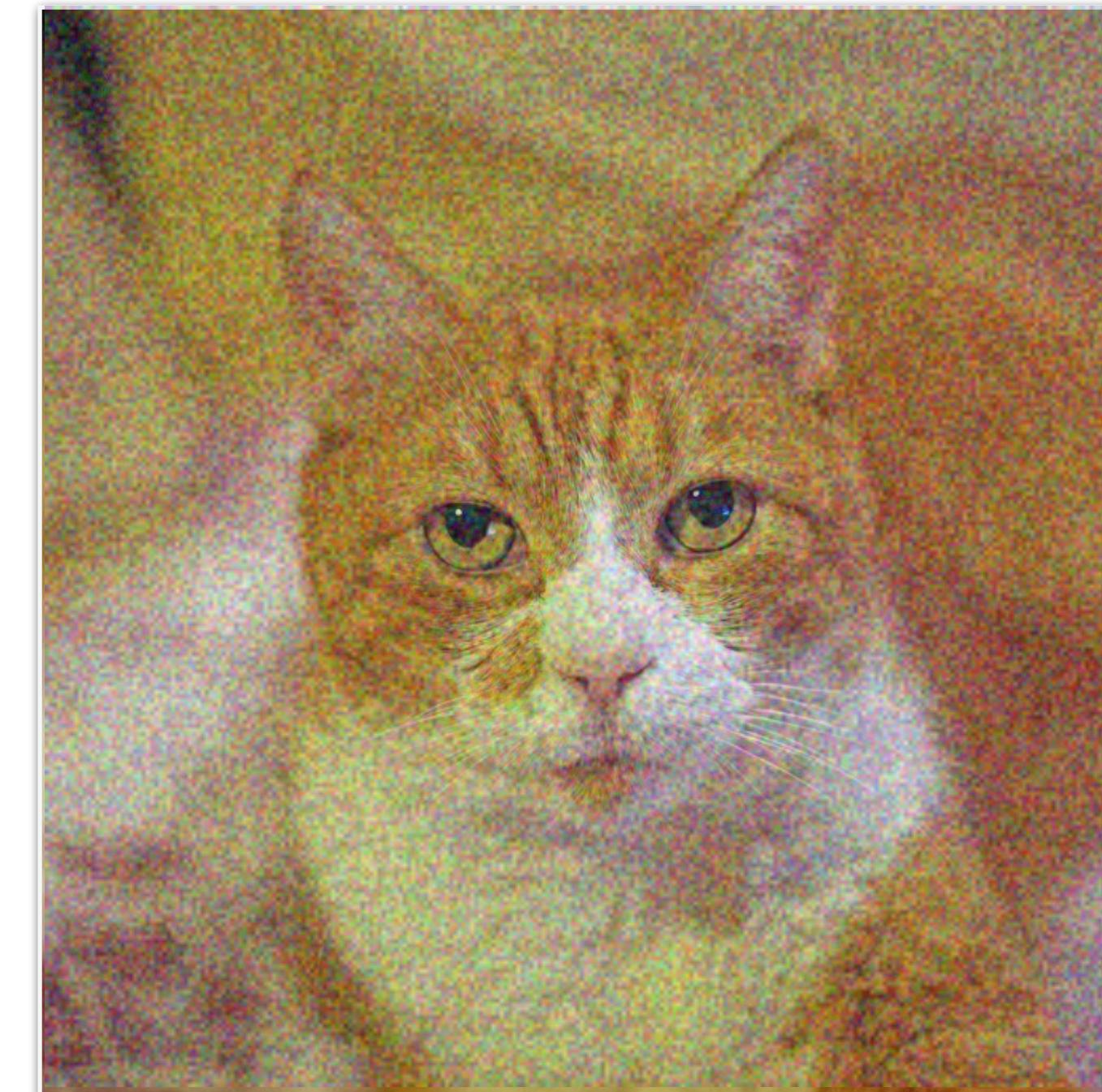
3. Optimize the image



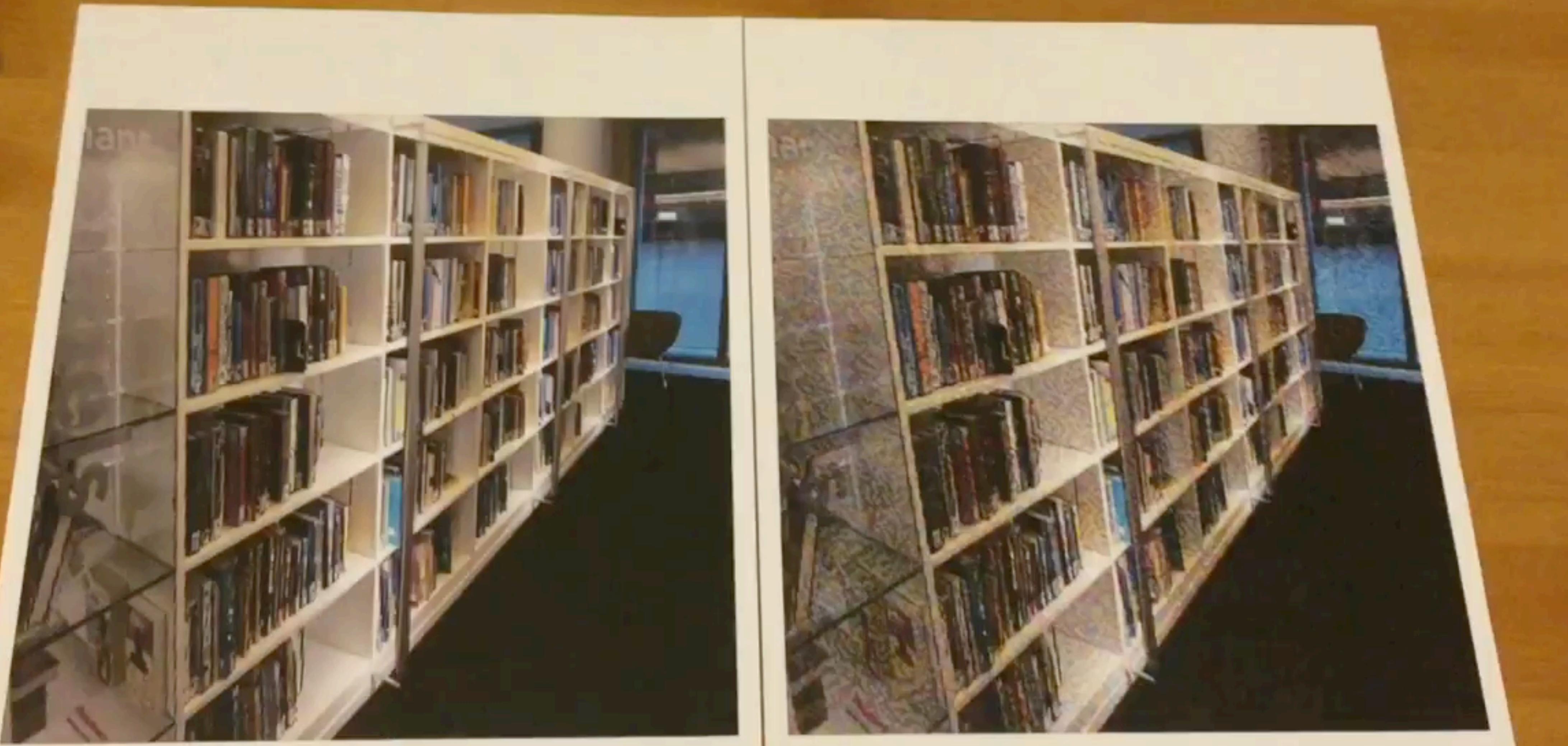
I. A. How to build adversarial examples and attack a network?



92% Cat



94% Iguana



Adversarial Examples In The Physical World

Kurakin A., Goodfellow I., Bengio S., 2016

I. B. How to defend against adversarial examples?

Types of attacks:

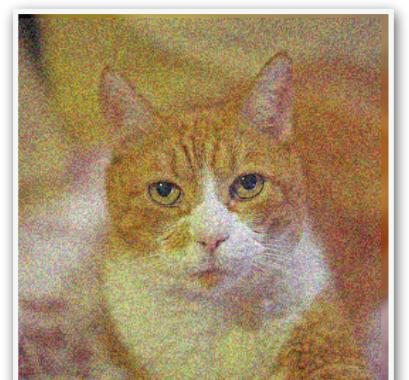
- Non-targeted attacks
- Targeted attacks

Knowledge of the attacker:

- White-box
- Black-box

Solution 1

- Train on correctly labelled adversarial examples



$x =$

$y = \text{cat}$

Solution 2

- Adversarial training $L_{new} = L(W, b, x, y) + \lambda L(W, b, x_{adv}, y)$

- Adversarial logit pairing $L_{new} = L(W, b, x, y) + \lambda \left\| f(x; W, b) - f(x_{adv}; W, b) \right\|_2^2$

Do neural networks actually understand the data?

II. Generative Adversarial Networks (GANs)

A. Motivation

B. G/D Game

C. Practical tips to train/evaluate GANs

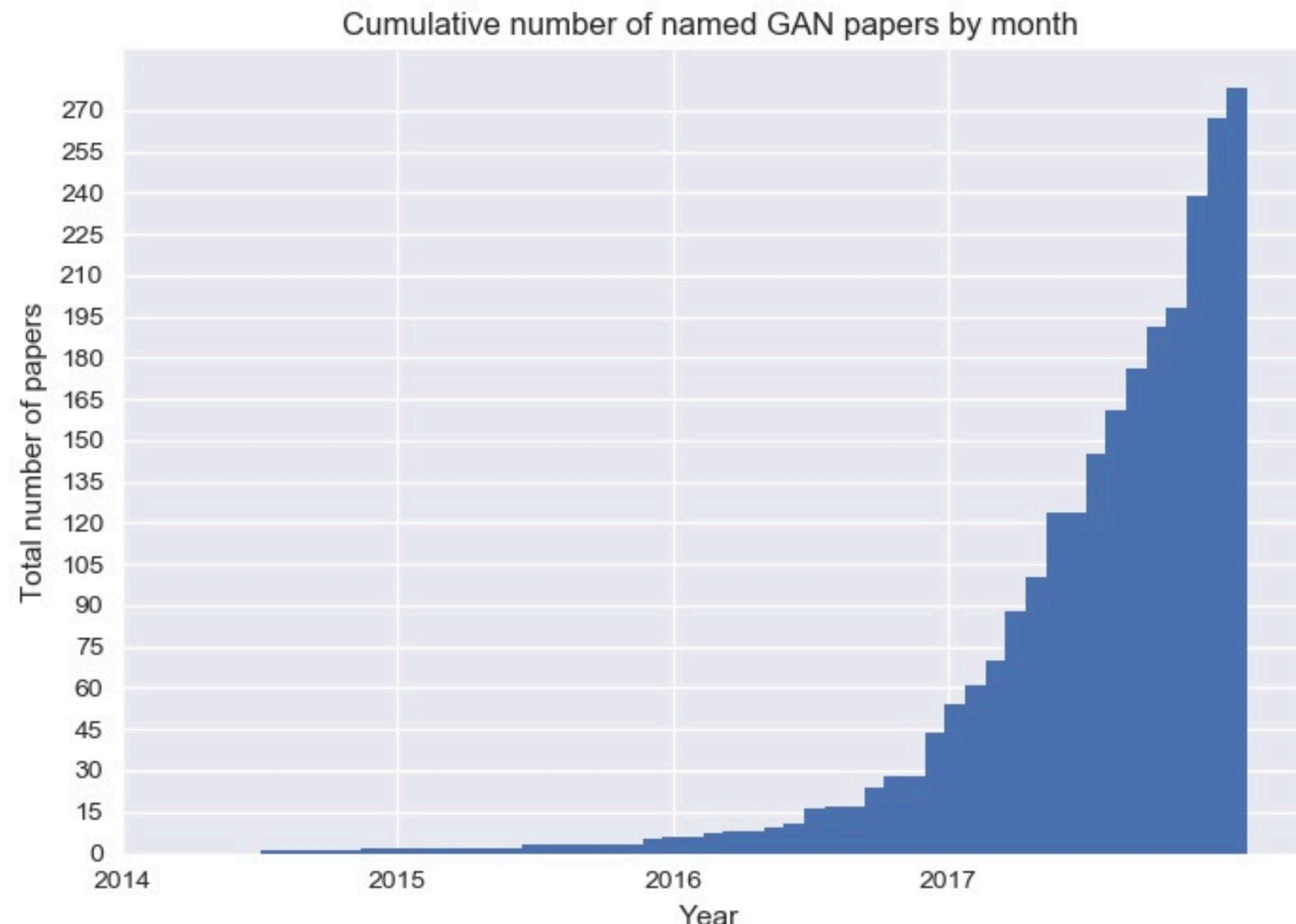
D. Interesting results

II.A - Motivation

Motivation: endowing computers with an understanding of our world.

Goal: collect a lot of data, use it to train a model to generate similar data from scratch.

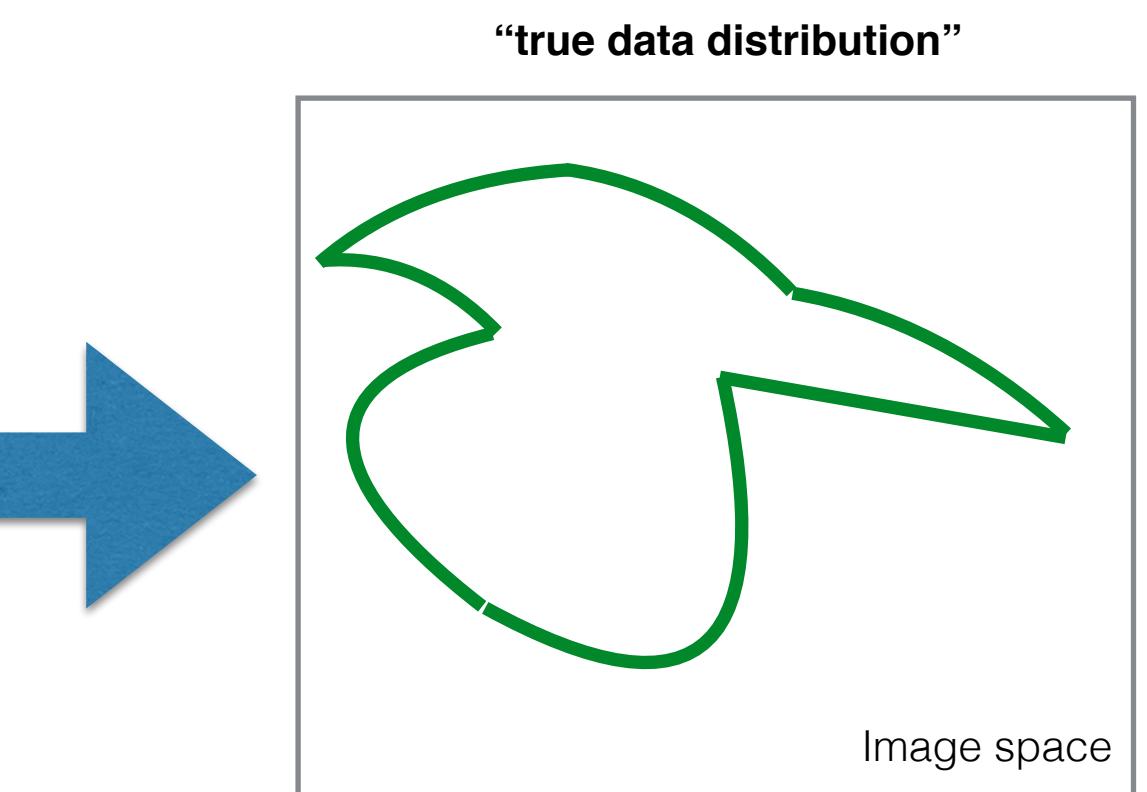
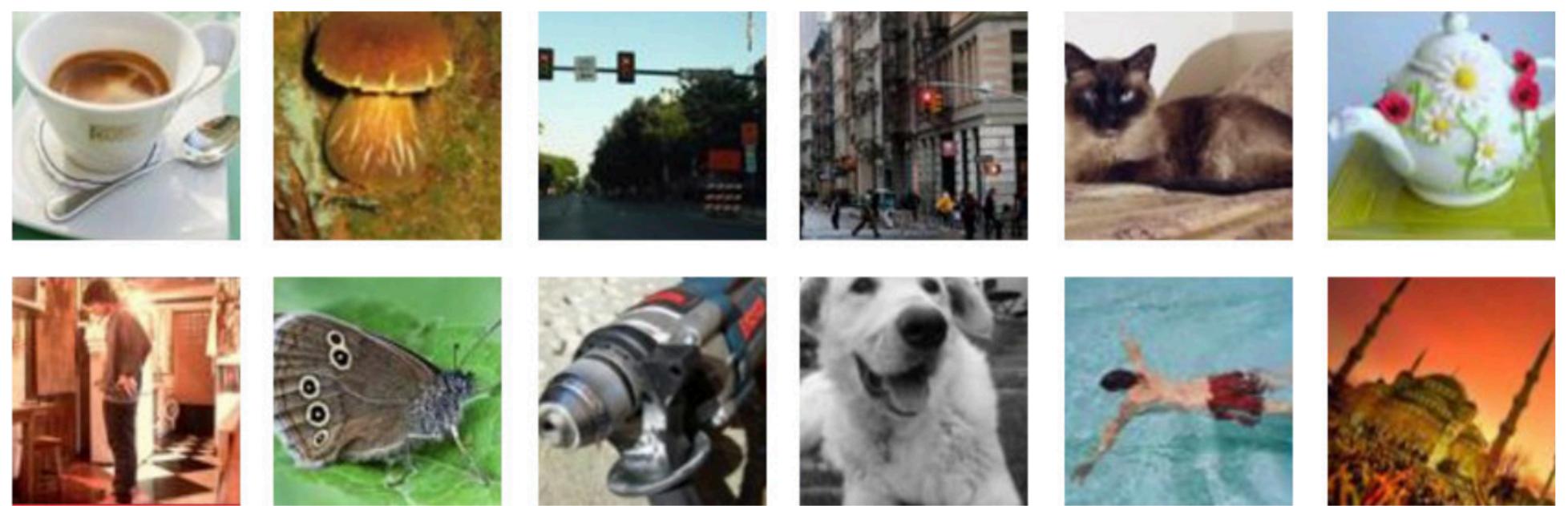
Intuition: number of parameters of the model << amount of data



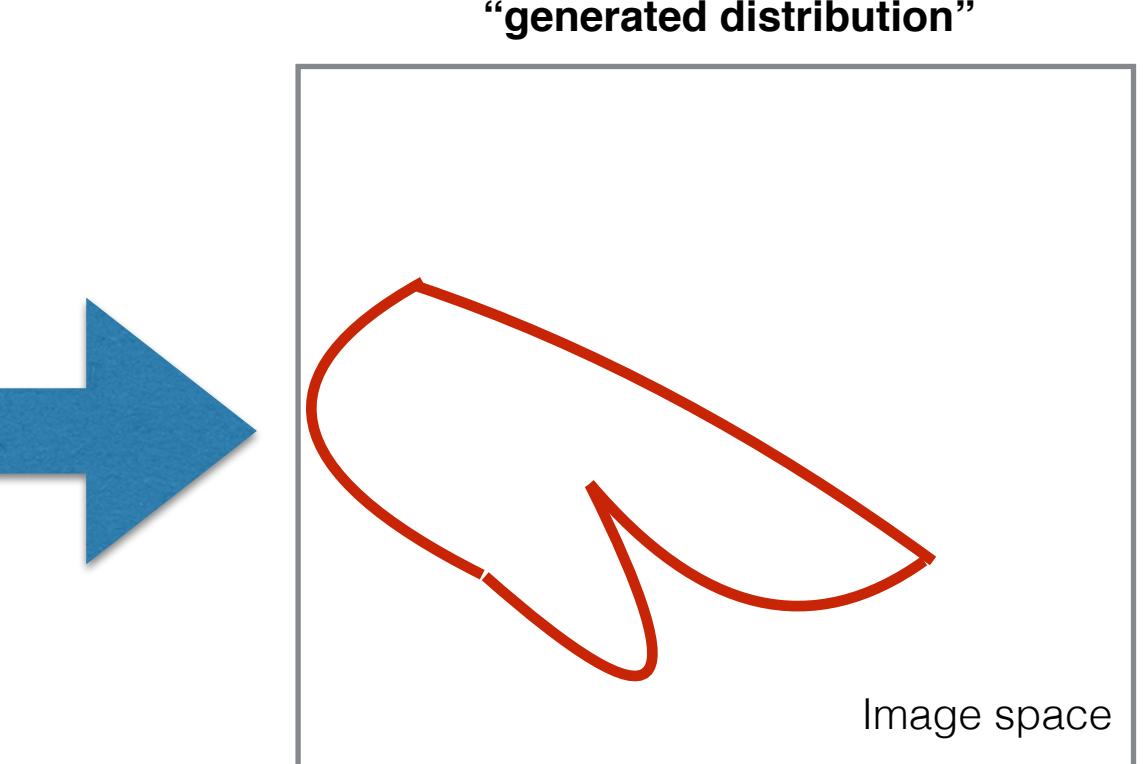
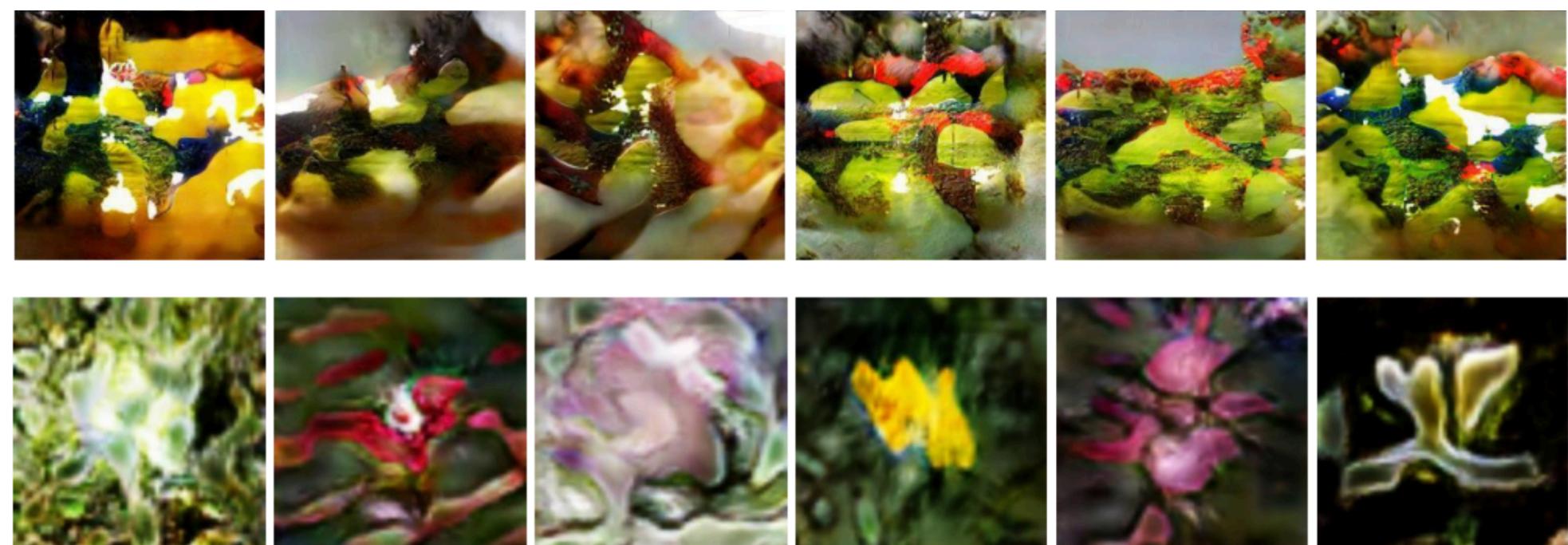
II.A - Motivation

Probability distributions:

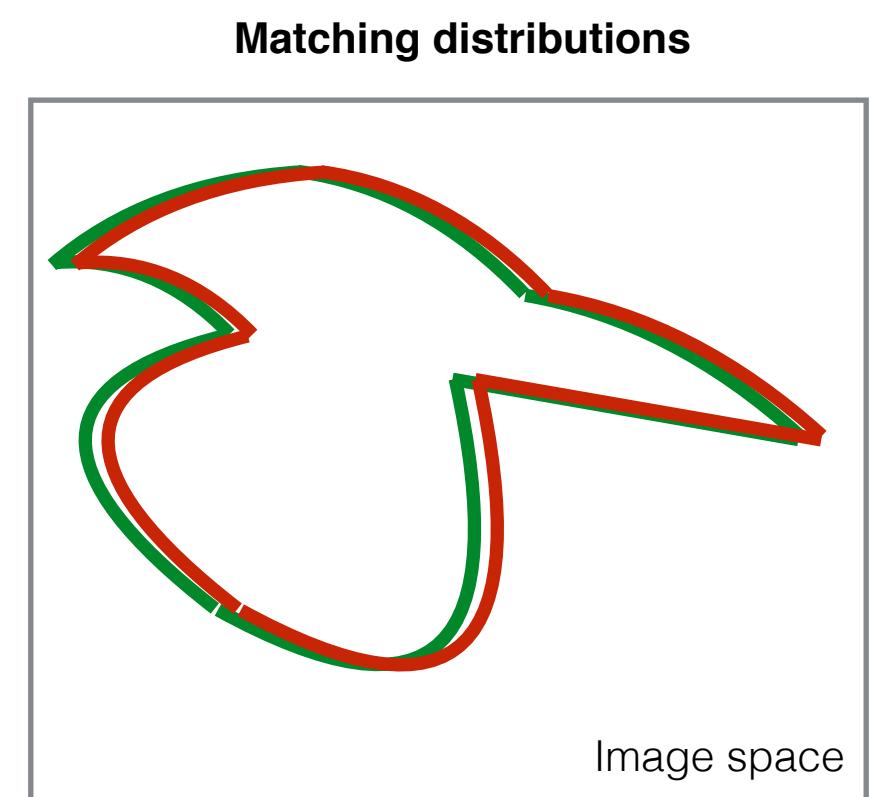
Samples from the “true data distribution”



Samples from the “generated distribution”



Goal

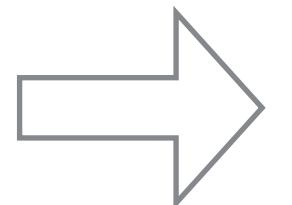
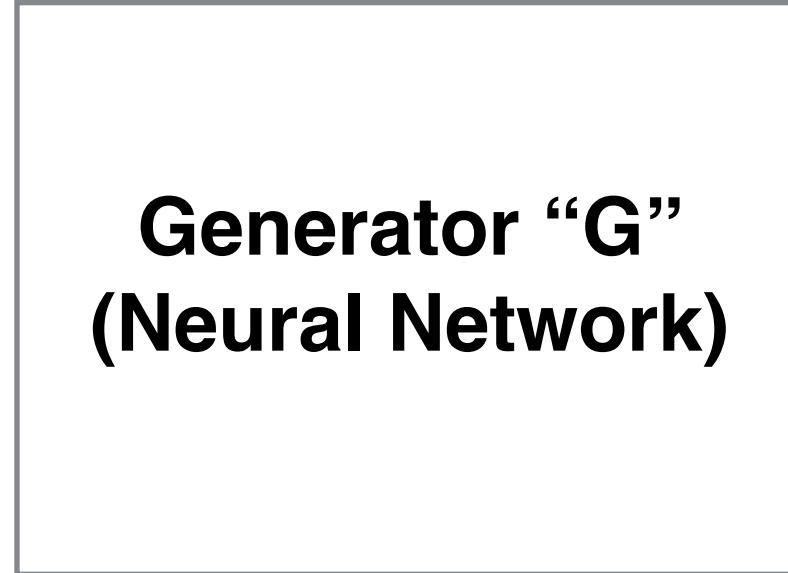


II.B - G/D Game

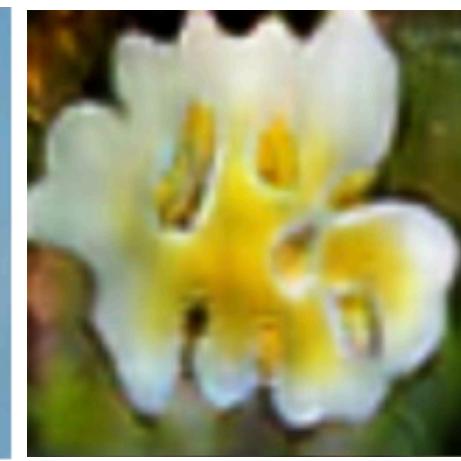
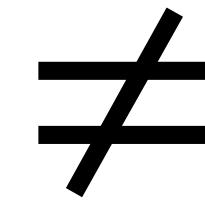
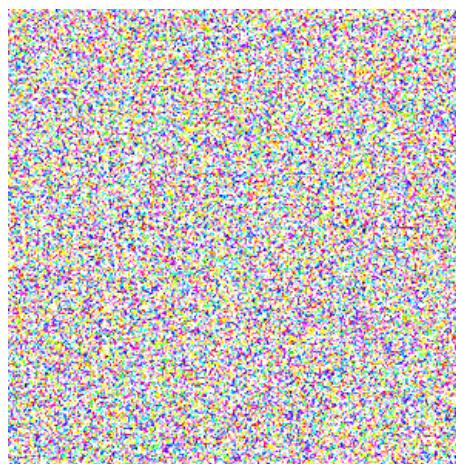
100-d
random code

$$\begin{pmatrix} 0.47 \\ \vdots \\ 0.19 \end{pmatrix}$$

z

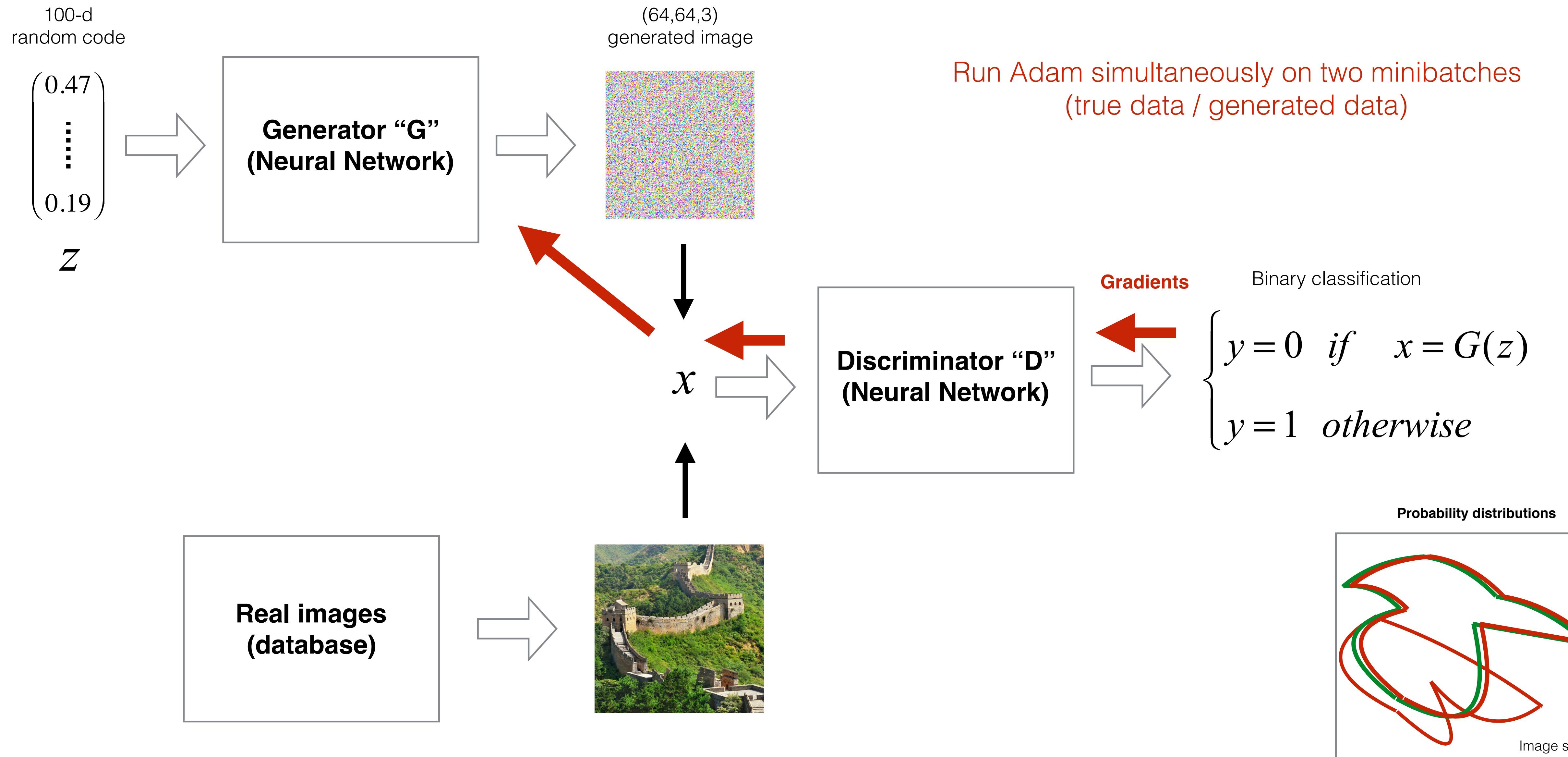


(64,64,3)
generated image

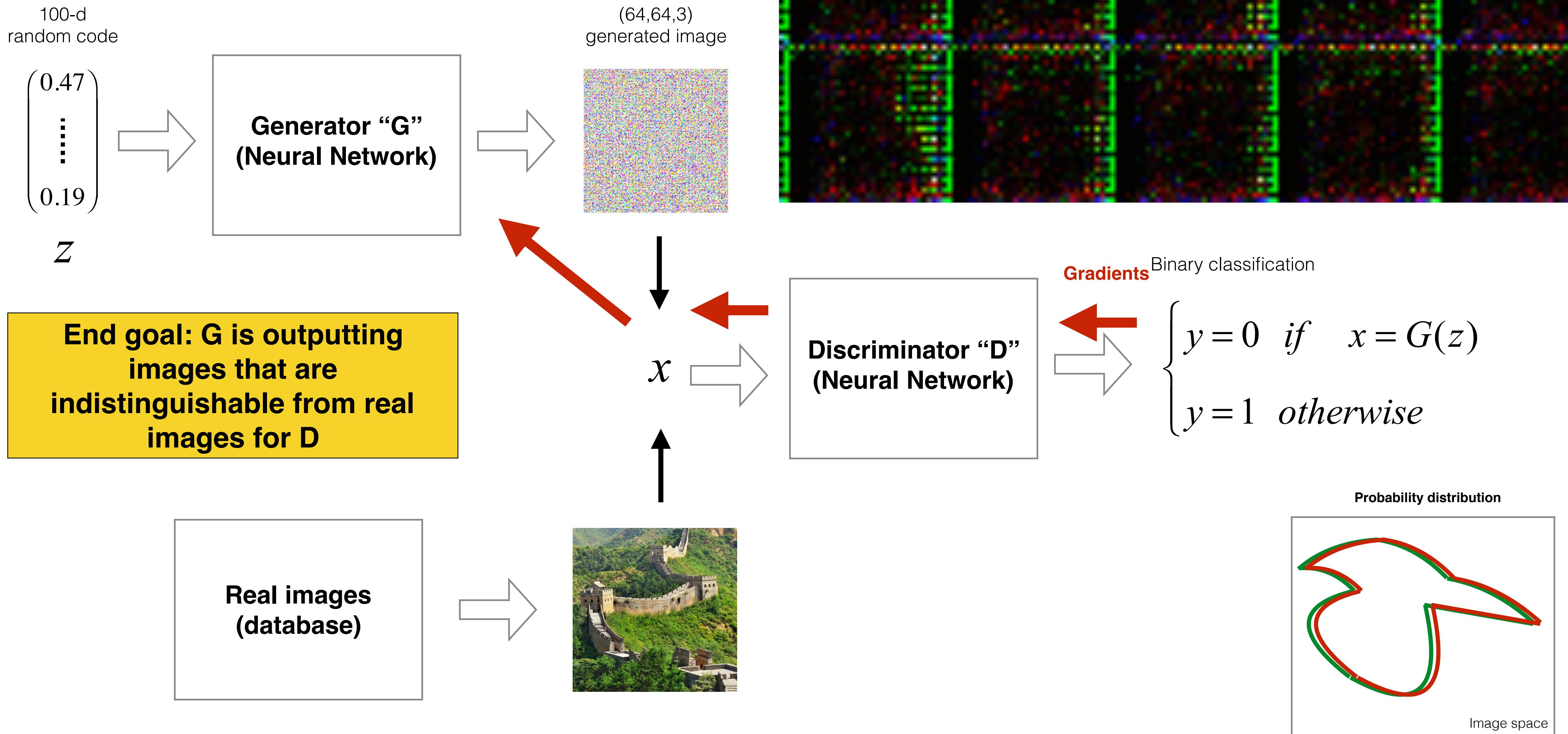


How can we train G to generate images from the true data distributions?

II.B - G/D Game



II.B - G/D



II.B - G/D Game

Training procedure, we want to minimize:

Labels: $\begin{cases} y_{real} & \text{is always 1} \\ y_{gen} & \text{is always 0} \end{cases}$

- The loss of the discriminator

$$J^{(D)} = \underbrace{-\frac{1}{m_{real}} \sum_{i=1}^{m_{real}} y_{real}^{(i)} \cdot \log(D(x^{(i)}))}_{\text{cross-entropy 1: "D should correctly label real data as 1"}} + \underbrace{-\frac{1}{m_{gen}} \sum_{i=1}^{m_{gen}} (1 - y_{gen}^{(i)}) \cdot \log(1 - D(G(z^{(i)})))}_{\text{cross-entropy 2: "D should correctly label generated data as 0"}}$$

- The loss of the generator

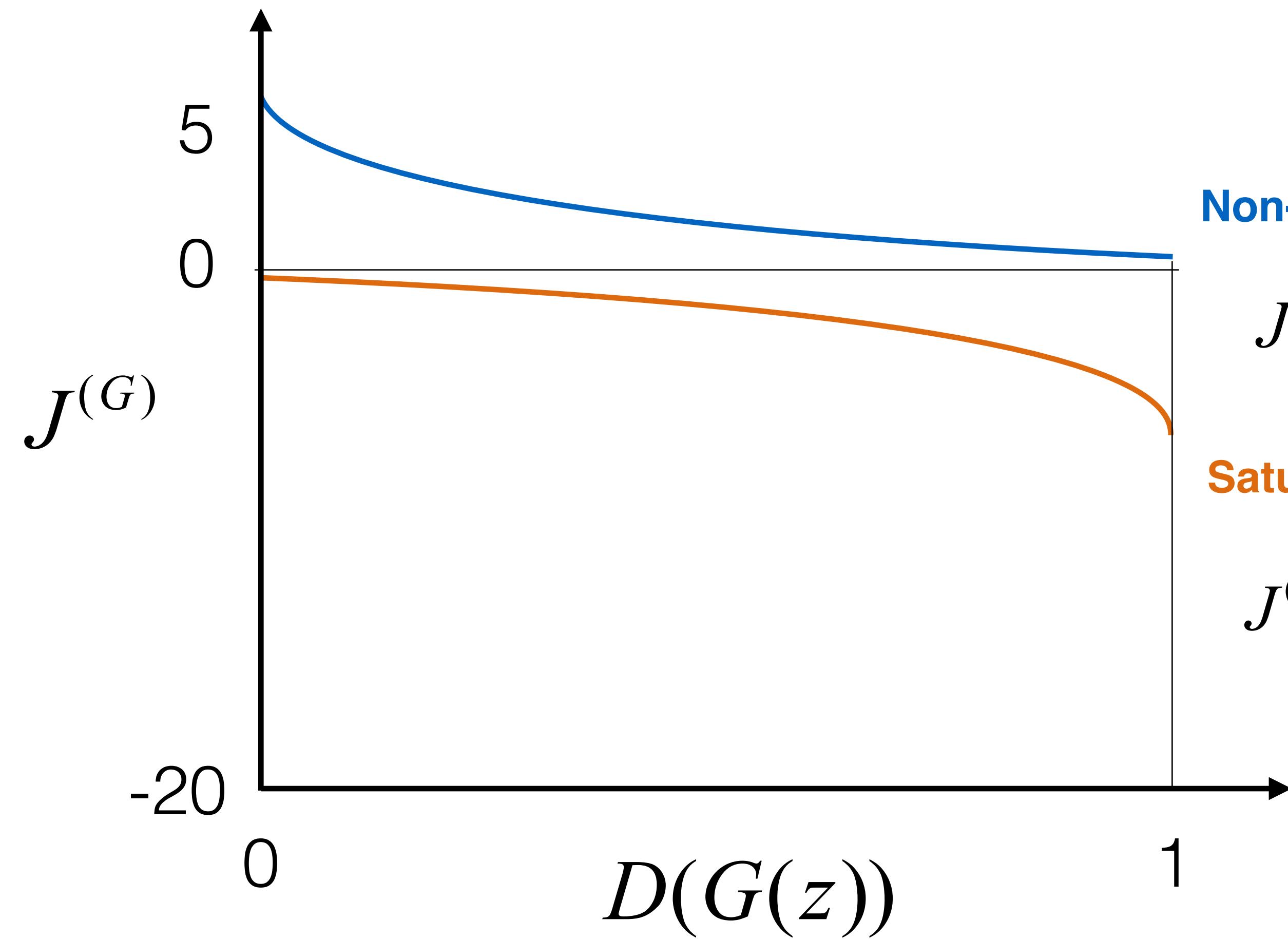
$$J^{(G)} = -J^{(D)} = \frac{1}{m_{gen}} \sum_{i=1}^{m_{gen}} \log(1 - D(G(z^{(i)})))$$

"G should try to fool D: by minimizing the opposite of what D is trying to minimize"

II.C - Training GANs

Saturating cost
for the generator:

$$\min \left[\frac{1}{m_{gen}} \sum_{i=1}^{m_{gen}} \log(1 - D(G(z^{(i)}))) \right] \Leftrightarrow \max \left[\frac{1}{m_{gen}} \sum_{i=1}^{m_{gen}} \log(D(G(z^{(i)}))) \right] \Leftrightarrow \min \left[-\frac{1}{m_{gen}} \sum_{i=1}^{m_{gen}} \log(D(G(z^{(i)}))) \right]$$



Non-saturating cost

$$J^{(G)} = -\frac{1}{m_{gen}} \sum_{i=1}^{m_{gen}} \log(D(G(z^{(i)})))$$

Saturating cost

$$J^{(G)} = \frac{1}{m_{gen}} \sum_{i=1}^{m_{gen}} \log(1 - D(G(z^{(i)})))$$

II.C - Training GANs

Note that:

$$\min \left[\frac{1}{m_{gen}} \sum_{i=1}^{m_{gen}} \log(1 - D(G(z^{(i)}))) \right] \Leftrightarrow \max \left[\frac{1}{m_{gen}} \sum_{i=1}^{m_{gen}} \log(D(G(z^{(i)}))) \right] \Leftrightarrow \min \left[-\frac{1}{m_{gen}} \sum_{i=1}^{m_{gen}} \log(D(G(z^{(i)}))) \right]$$

New training procedure, we want to minimize:

$$J^{(D)} = \underbrace{-\frac{1}{m_{real}} \sum_{i=1}^{m_{real}} y_{real}^{(i)} \cdot \log(D(x^{(i)}))}_{\text{cross-entropy 1: "D should correctly label real data as 1"} } - \underbrace{\frac{1}{m_{gen}} \sum_{i=1}^{m_{gen}} (1 - y_{gen}^{(i)}) \cdot \log(1 - D(G(z^{(i)})))}_{\text{cross-entropy 2: "D should correctly label generated data as 0"} }$$

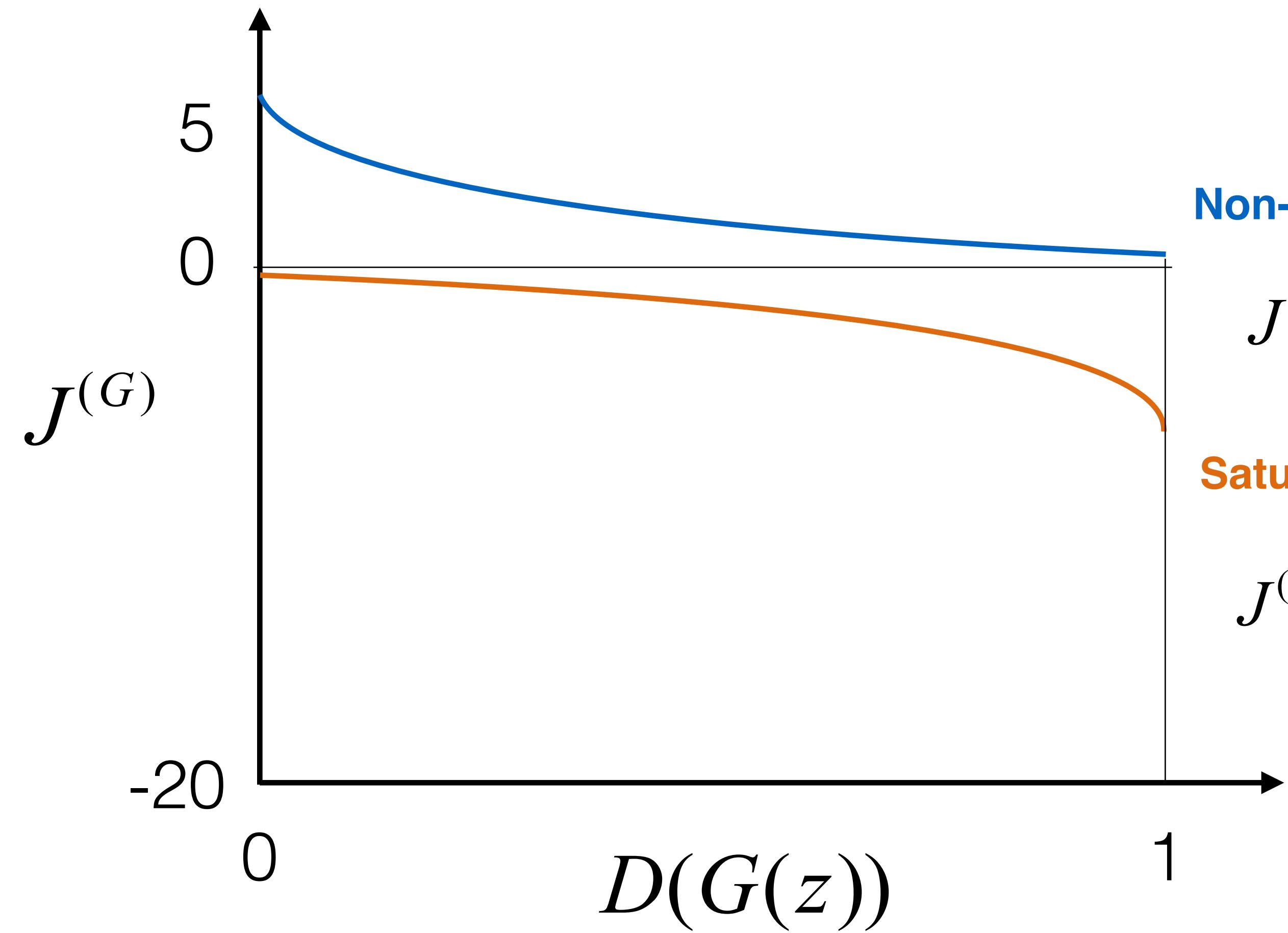
$$J^{(G)} = -\frac{1}{m_g} \sum_{i=1}^{m_g} \log(D(G(z^{(i)})))$$

“G should try to fool D: by minimizing this”

II.C - Training GANs

Simultaneously training G/D?

```
for num_iterations:  
    for k iterations:  
        update D  
        update G
```



Non-saturating cost

$$J^{(G)} = -\frac{1}{m_g} \sum_{i=1}^{m_g} \log(D(G(z^{(i)})))$$

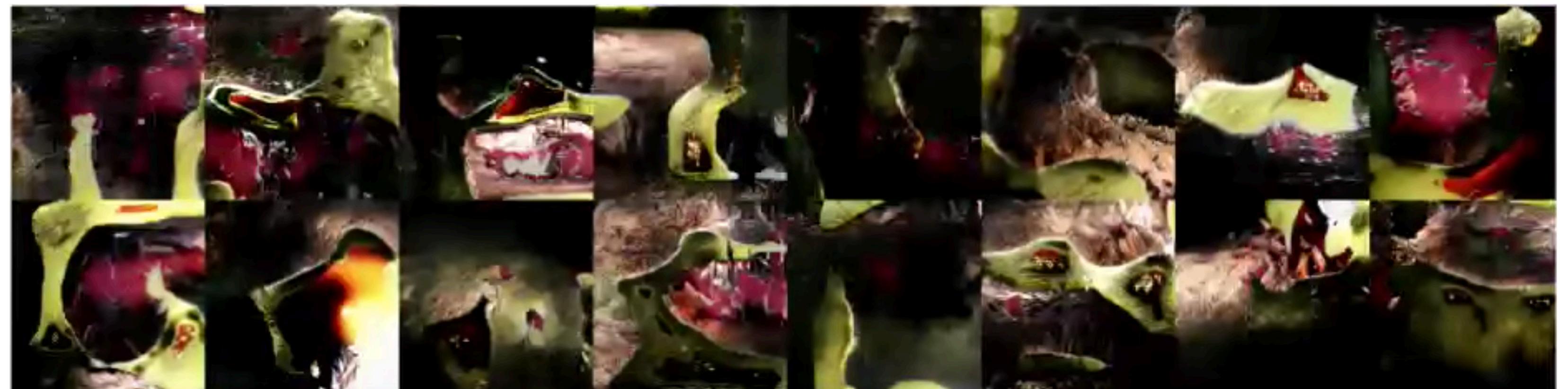
Saturating cost

$$J^{(G)} = \frac{1}{m_g} \sum_{i=1}^{m_g} \log(1 - D(G(z^{(i)})))$$

II.C - Training GANs

BatchNorm with GANs:

Generated images
(batch 1)



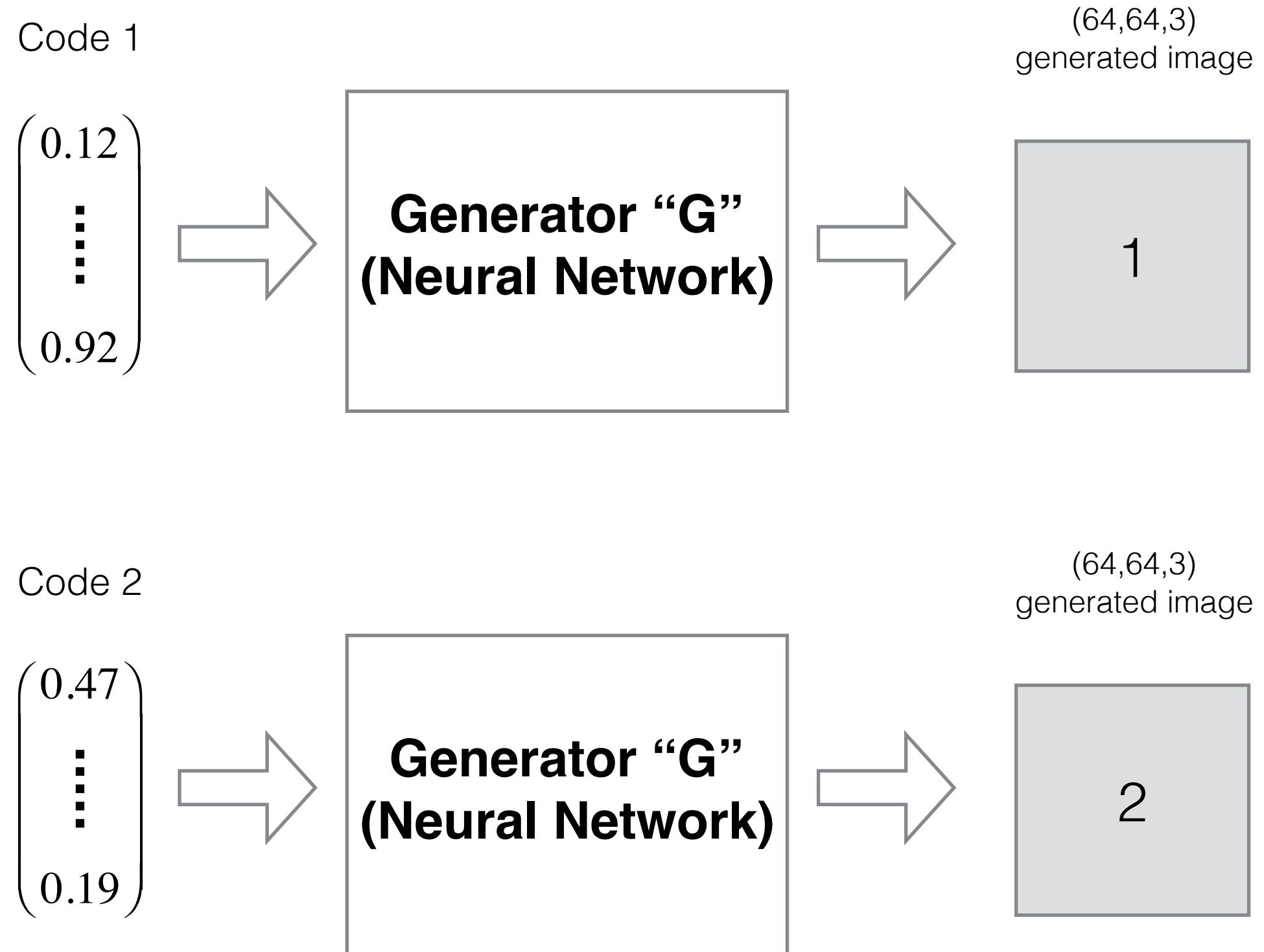
Generated images
(batch 2)



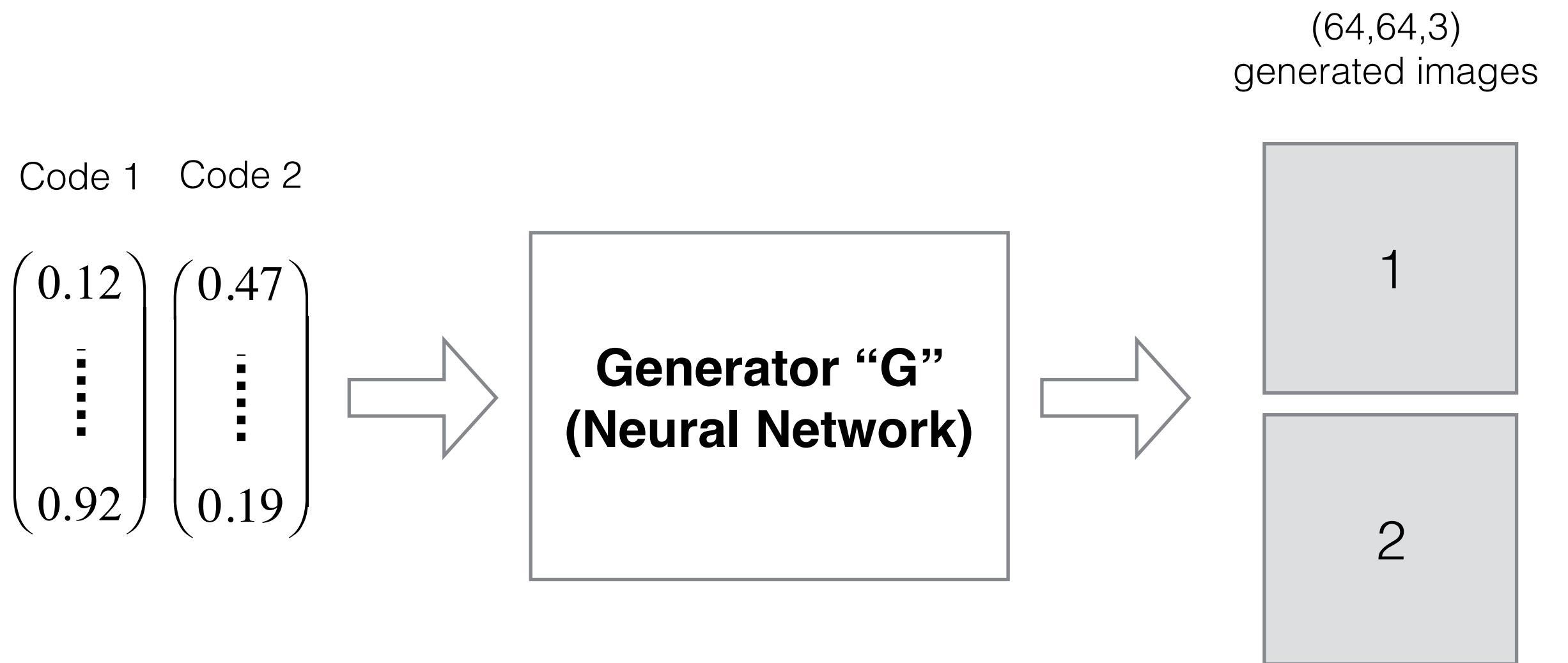
II.C - Training GANs

BatchNorm with GANs:

Assume no batchnorm



Assume batchnorm



II.C - Training GANs

BatchNorm with GANs:

BatchNorm

$$Z = \{z^{(1)}, \dots, z^{(m)}\}$$

$$\mu_B = \frac{1}{m} \sum_{i=1}^m z^{(i)}$$

$$\sigma_B^2 = \frac{1}{m} \sum_{i=1}^m (z^{(i)} - \mu_B)^2$$

$$z_{norm}^{(i)} = \frac{z^{(i)} - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$

$$\tilde{z}^{(i)} = \gamma z_{norm}^{(i)} + \beta$$

Reference BatchNorm

$$R = \{r^{(1)}, \dots, r^{(m)}\}$$

$$Z = \{z^{(1)}, \dots, z^{(m)}\}$$

$$\mu_B = \frac{1}{m} \sum_{i=1}^m r^{(i)}$$

$$\sigma_B^2 = \frac{1}{m} \sum_{i=1}^m (r^{(i)} - \mu_B)^2$$

$$z_{norm}^{(i)} = \frac{z^{(i)} - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$

$$\tilde{z}^{(i)} = \gamma z_{norm}^{(i)} + \beta$$

Virtual BatchNorm

$$R = \{r^{(1)}, \dots, r^{(m)}\}$$

$$Z = \{z^{(1)}, \dots, z^{(m)}\}$$

For k = 1....m

$$\mu_B = \frac{1}{m+1} \left(z^{(k)} + \sum_{i=1}^m r^{(i)} \right)$$

$$\sigma_B^2 = \frac{1}{m+1} \left((z^{(k)} - \mu_B)^2 + \sum_{i=1}^m (r^{(i)} - \mu_B)^2 \right)$$

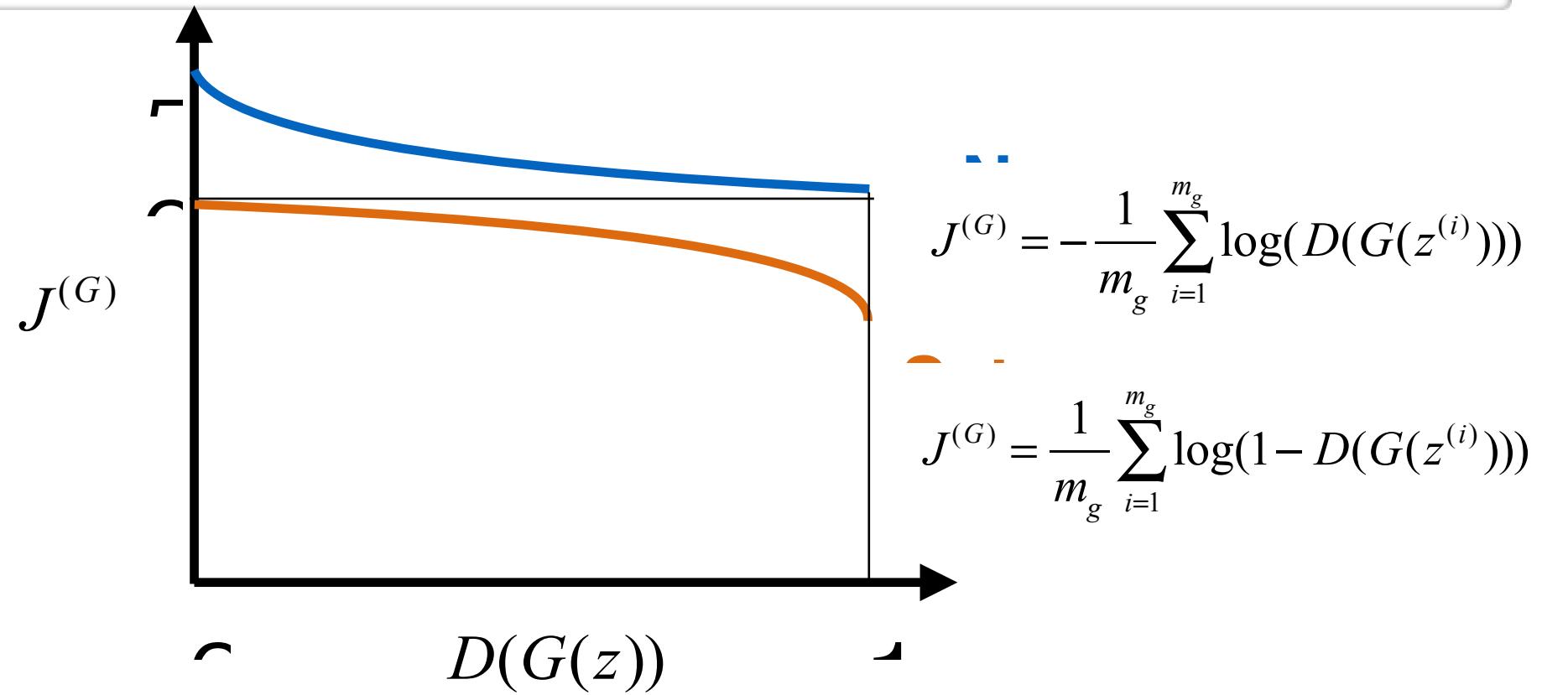
$$z_{norm}^{(k)} = \frac{z^{(k)} - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$

$$\tilde{z}^{(k)} = \gamma z_{norm}^{(k)} + \beta$$

II.C - Training GANs

Recap: GANs' training tips

- Use the non-saturated cost function
- Keep D up-to-date with respect to G (k update for D / 1 update for G)
- Use Virtual Batchnorm
- (not presented but important) One-sided label smoothing



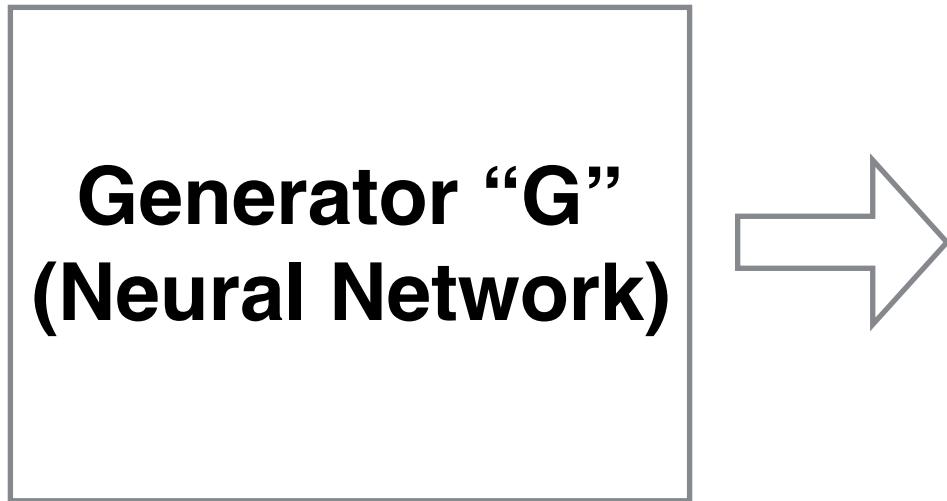
And a lot more, GANs are hard to train!

II.D - Interesting results

Operation on codes

Code 1

$$\begin{pmatrix} 0.12 \\ \vdots \\ 0.92 \end{pmatrix}$$

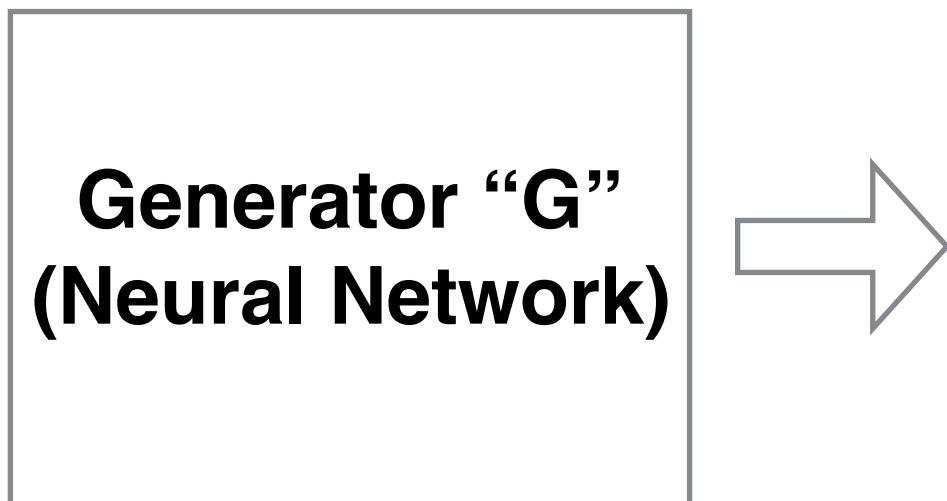


(64,64,3)
generated image



Code 2

$$\begin{pmatrix} 0.47 \\ \vdots \\ 0.19 \end{pmatrix}$$

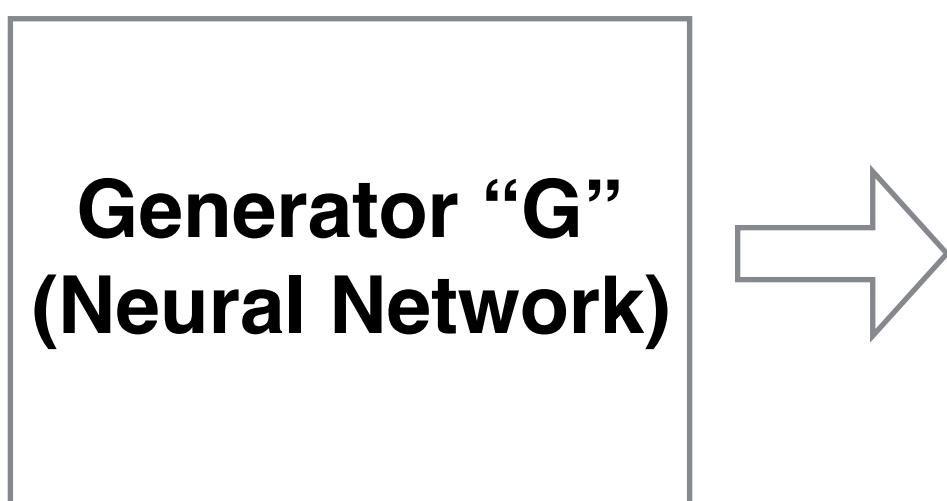


(64,64,3)
generated image

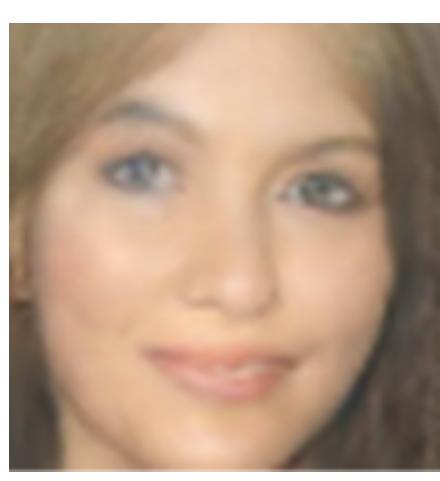


Code 3

$$\begin{pmatrix} 0.42 \\ \vdots \\ 0.07 \end{pmatrix}$$

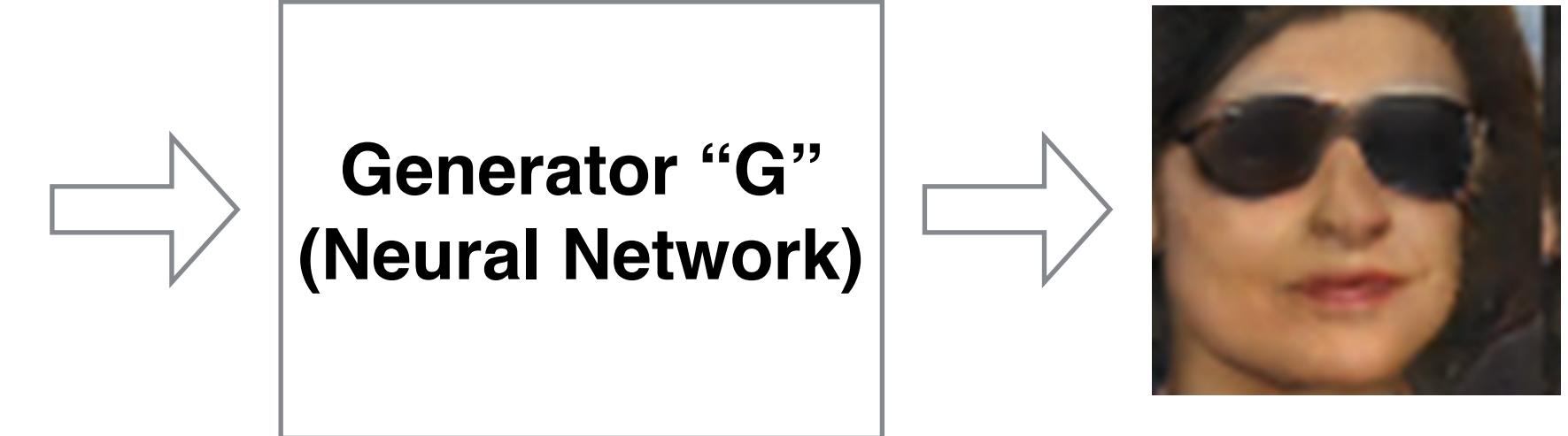


(64,64,3)
generated image



Code 1 Code 2 Code 3

$$\begin{pmatrix} 0.12 \\ \vdots \\ 0.92 \end{pmatrix} - \begin{pmatrix} 0.47 \\ \vdots \\ 0.19 \end{pmatrix} + \begin{pmatrix} 0.42 \\ \vdots \\ 0.07 \end{pmatrix}$$

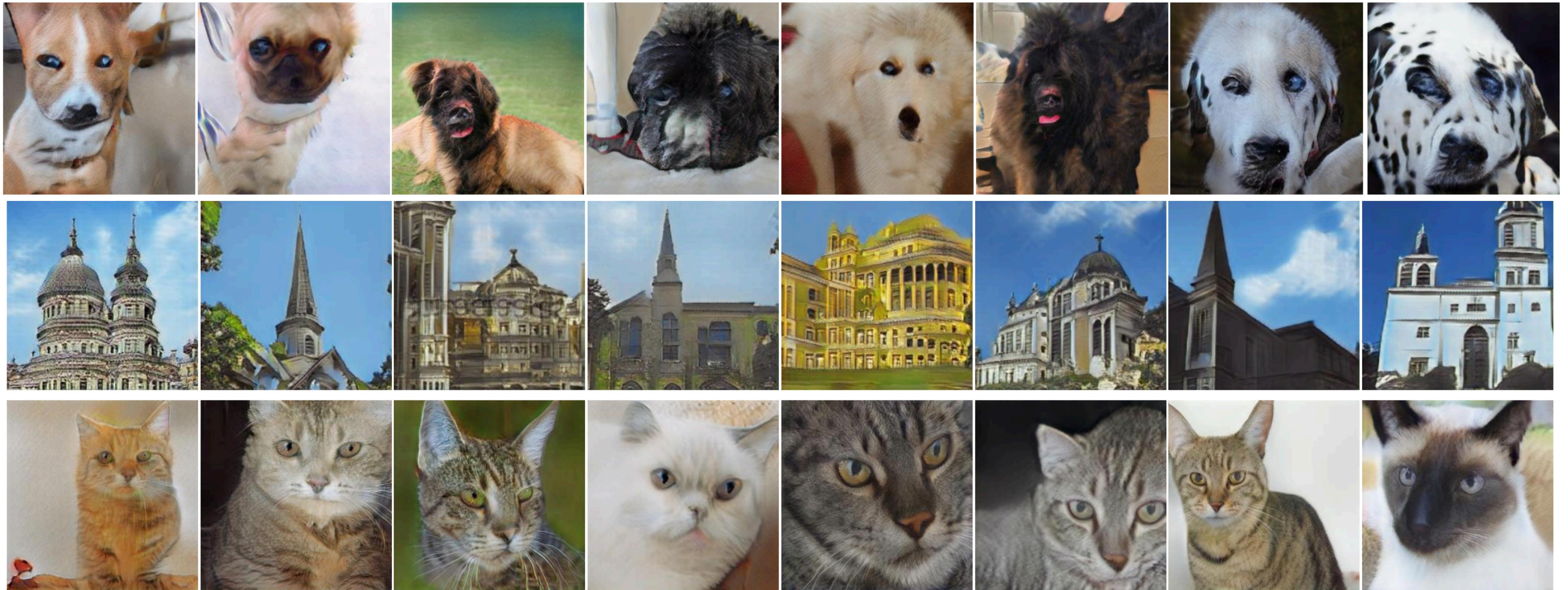


Man with glasses - man + woman = woman with glasses

II.D - Interesting results

Image Generation:

Samples from the “generated distribution”



II.D - Interesting results

Pix2Pix:

<https://affinelayer.com/pixsrv/>

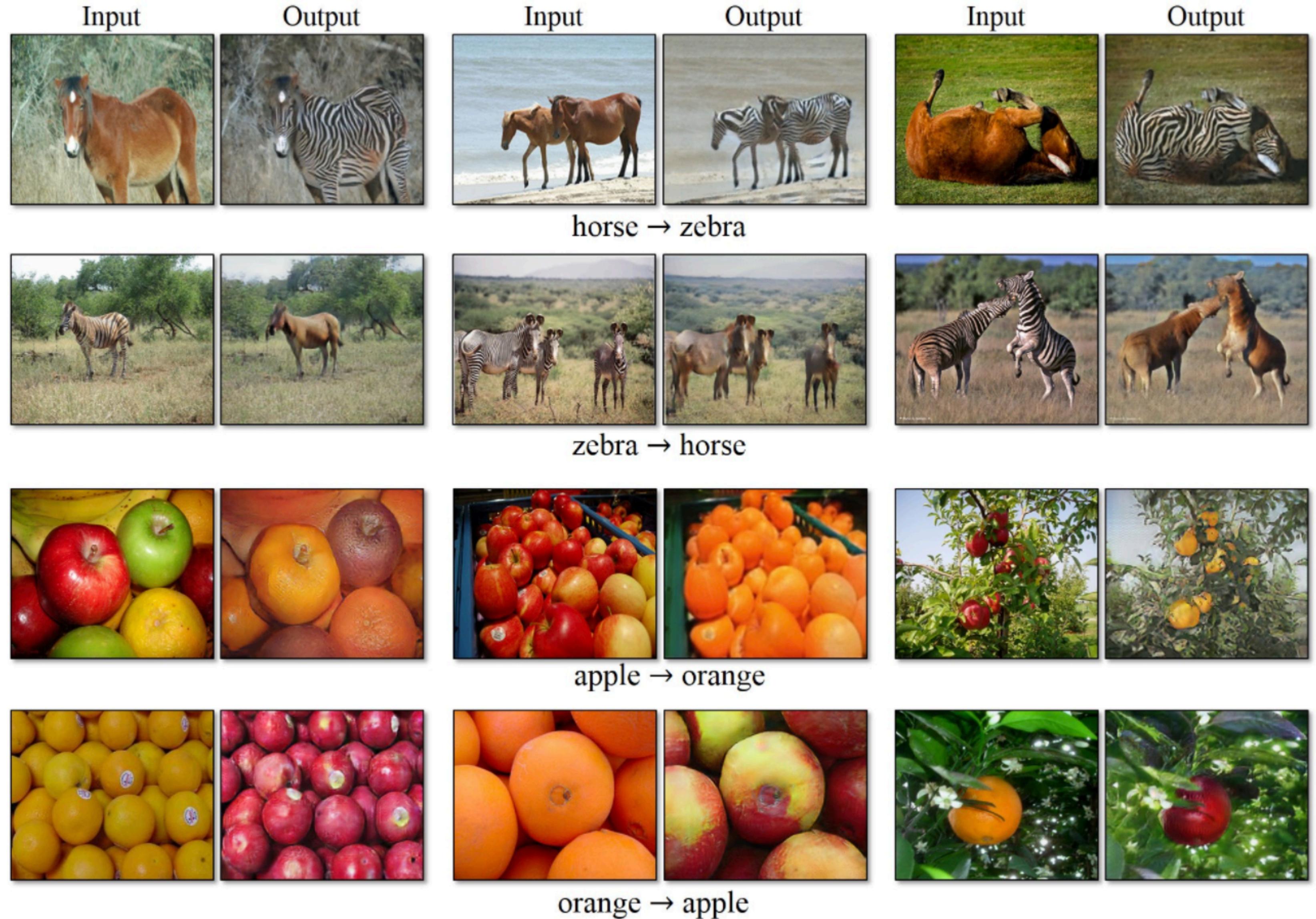
II.D - Interesting results

Super-resolution image:



II.D - Interesting results

CycleGANs:

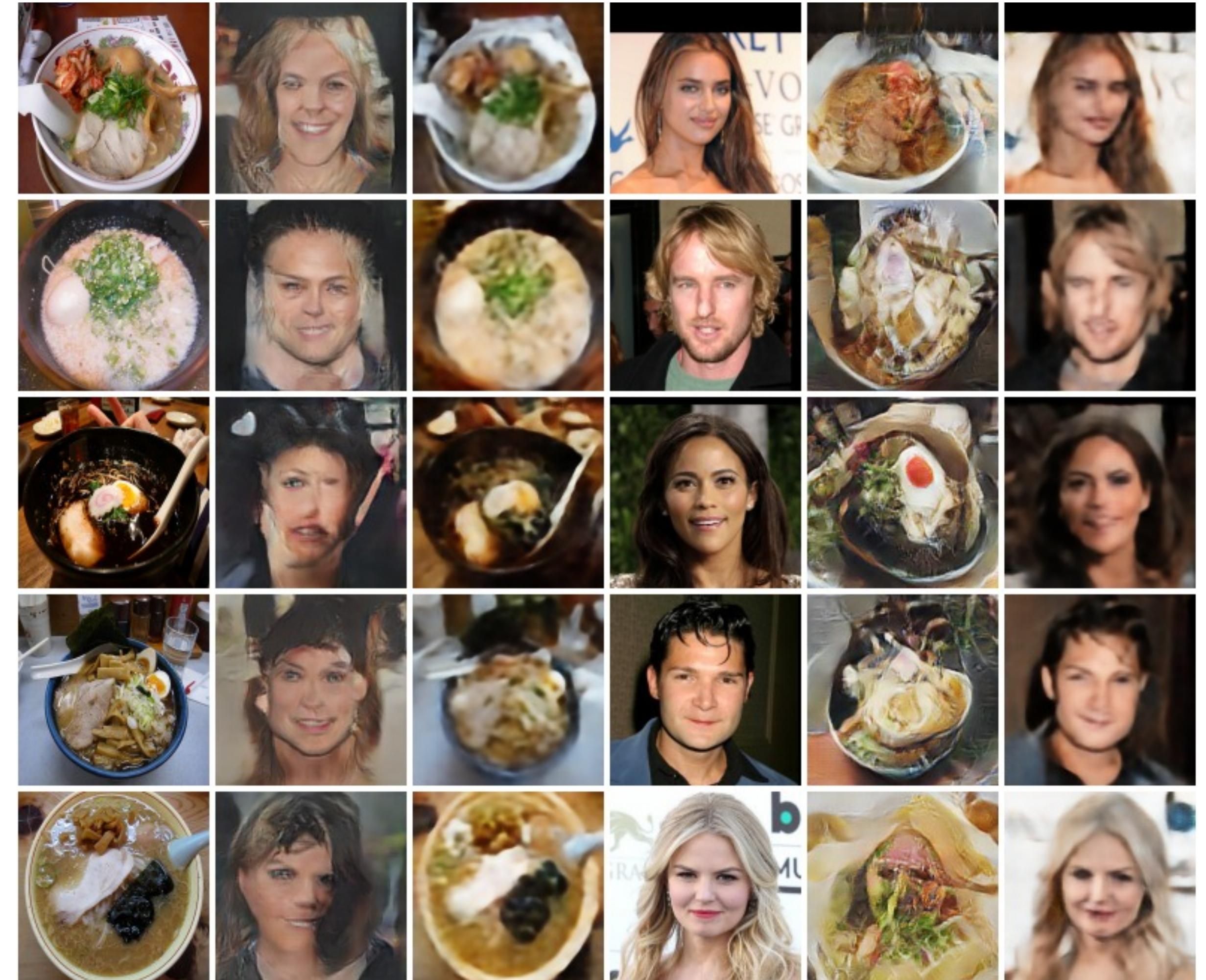


II.D - Interesting results

CycleGANs:

<https://hardikbansal.github.io/CycleGANBlog/>

Face2ramen



Announcements

For Tuesday 05/01, 9am:

C2M3

- Quiz: Hyperparameter tuning, Batch Normalization, Programming Frameworks
- Programming assignment: Tensorflow

C3M1 and C3M2

- Quiz: Bird recognition in the city of Peacetopia (case study)
- Quiz: Autonomous driving (case study)

For Friday 02/16, 9am:

- Hands-on session this Friday

Check out the project
example code!
(cs230-stanford.github.io)

Meet with your mentor (TA), you'll receive a Calendly invite.