

**WHAT IF WE TRAIN  
THAT 2-ND NETWORK**



**TO HELP US TRAIN  
THE FIRST NETWORK**

imgflip.com

Lets train a  
model to  
tell us  
whether a  
generated  
image is  
good  
enough or  
not.

# Generative Adversarial Networks

This is the high level overview.  
GENERATOR VS  
DISCRIMINATOR.

Generator



Generate image  
(should be plausible)

content

feedback

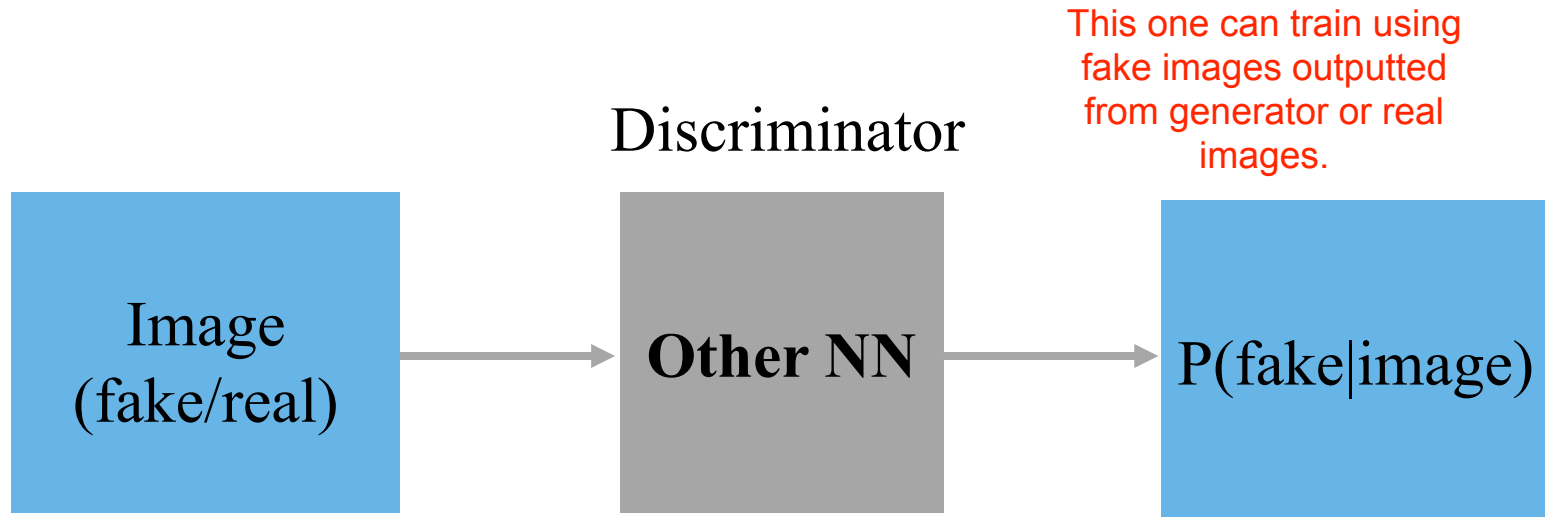
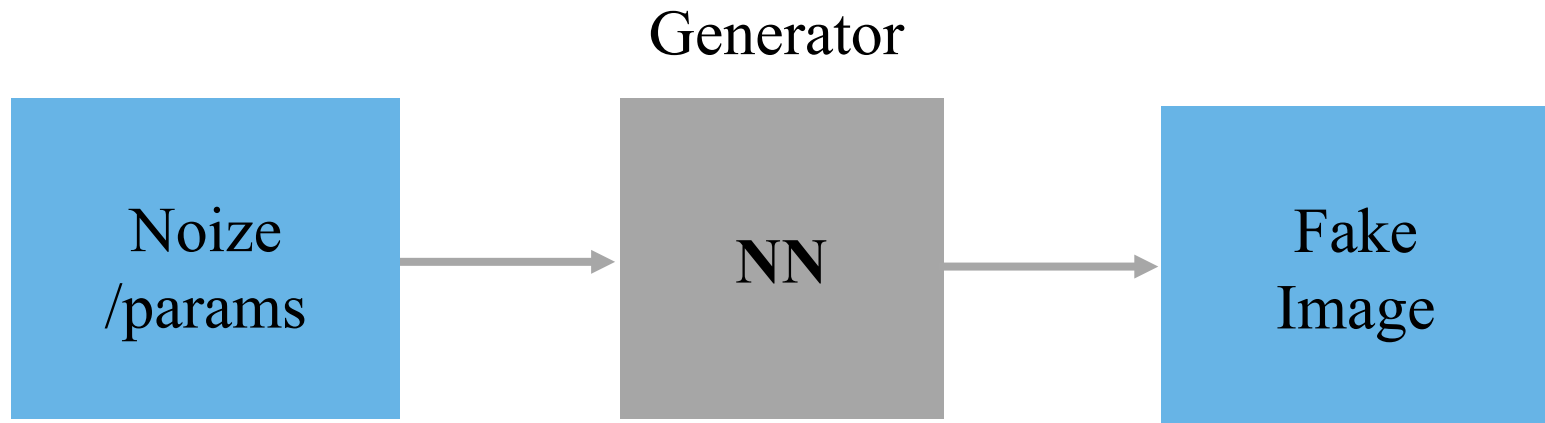
Discriminator



Tell if image is plausible  
(image)  $\rightarrow$  P(fake)

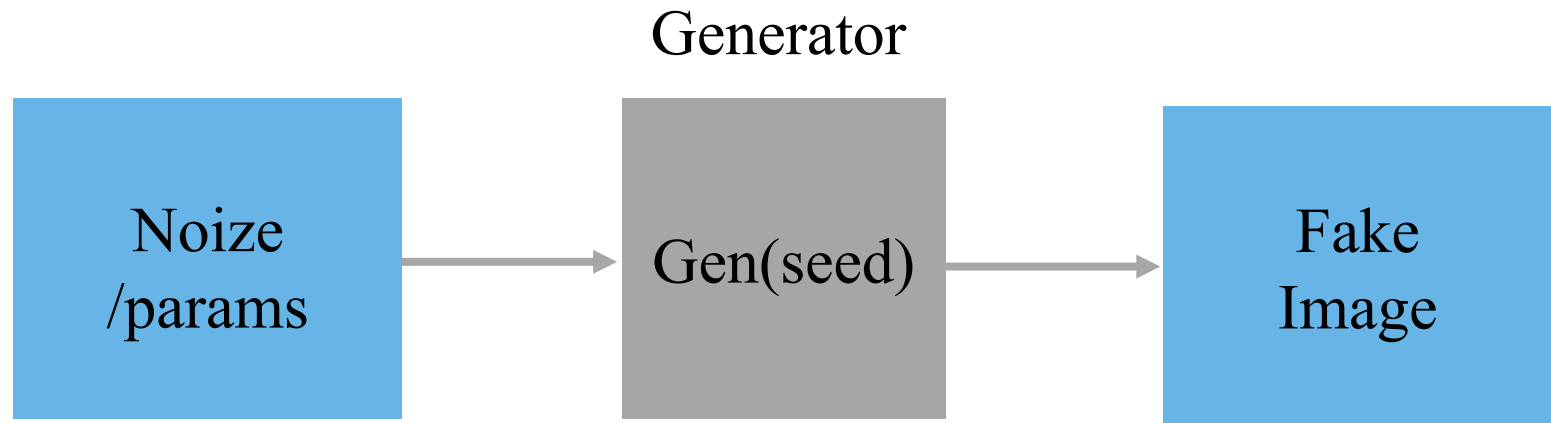
called discriminator  
because it discriminates  
between real images and  
the generated images.

# Generative Adversarial Networks

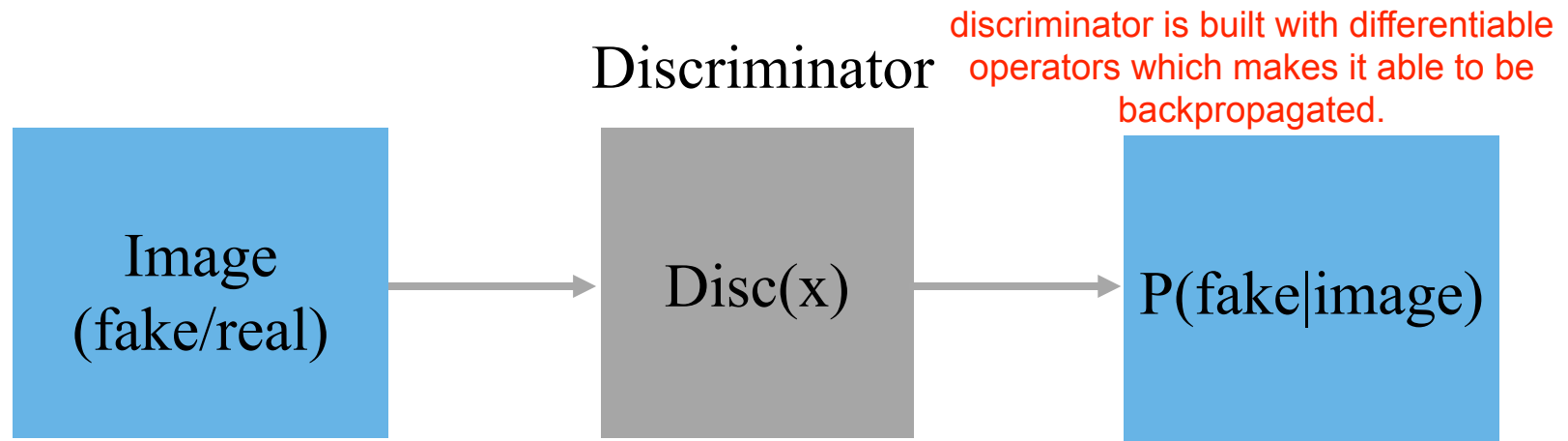


The discriminator can train on positive examples from all real images, and negative examples using all generated fake images.

# Generative Adversarial Networks



$$L_G = -\log[1 - \text{Disc}(\text{Gen}(\text{seed}))]$$



$$L_D = -\log[1 - \text{Disc}(\text{realdata})] - \log \text{Disc}(\text{Gen}(\text{seed}))$$

# Generative Adversarial Networks

## Algorithm

- sample noise  $\mathbf{z}$  and images  $\mathbf{x}$
- for  $k$  in  $1 \dots K$ 
  - Train discriminator( $\mathbf{x}$ ), discriminator(generator( $\mathbf{z}$ ))

- For  $m$  in  $1 \dots M$

<sup>^these are the  
positive samples.</sup>

<sup>^these are the negative samples.</sup>

- Train generator( $\mathbf{z}$ )

training steps:

Generative models are unstable. You have two models that 'hate each other'. And if one of them wins, you have to start the process all over again.

If discriminator wins (can train faster than generator), then the gradients vanish - the sigmoids that are used to compute  $P(\text{real} | \text{image})$  are really close to 1 or 0.

Thus, they have very small gradient.

If generator wins (constantly train faster than discriminator), then it can start learning the wrong things.

The generator can learn non-sensical stuff.

(0). initialize generator and discriminator weights randomly.

(1). Train discriminator to classify actual images against images generated by (untrained) generator

(2). Train generator to generate images that fool discriminator into believing they're real.

repeat (1) and then (2) again.

This cycle continues. This is why it's called ADVERSARIAL. it is as if they are competing against each other.

# Generative Adversarial Networks

