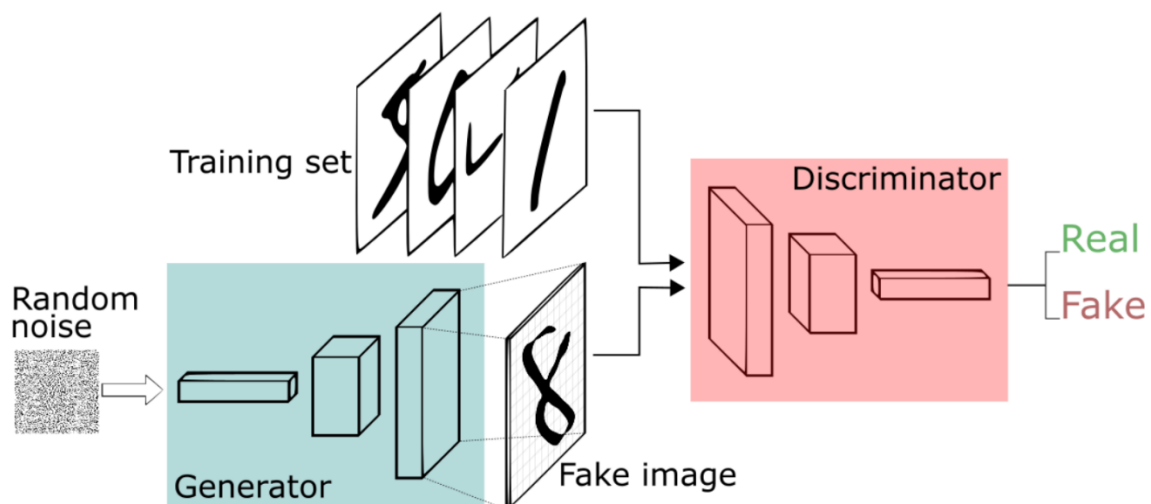# Independent Study IMT 600

## *Generative Adversarial Networks*

Generative adversarial networks (GANs) represents one of the most exciting recent innovation in deep learning. They are used widely in image generation, video generation and voice generation.

A GAN is a generative model in which two neural networks are competing in a typical game theory scenario. The first neural network is the **generator**, responsible of generating new synthetic data instances that resemble your training data, while its adversary, the **discriminator** tries to distinguish between real (training) and fake (artificially generated) samples generated by the generator. The mission of the generator is to try fooling the discriminator, and the discriminator tries to resist from being fooled. That's why the system as a whole is described as **adversarial**.

As shown in the figure below, the generator's input is simply a random noise, while, only the discriminator has access to the training data for classification purposes. The generator keeps improving its output based, exclusively, on the feedback of the discriminator network (positive in case of match with training data and negative if there is no match).



Generative Adversarial Network Architecture

**How GANs work**

One neural network, called the *generator*, generates new data instances, while the other, the *discriminator*, evaluates them for authenticity; i.e. the discriminator decides whether each instance of data that it reviews belongs to the actual training dataset or not.

Let's say we're trying to do something more banal than mimic the Mona Lisa. We're going to generate hand-written numerals like those found in the MNIST dataset, which is taken from the real world. The goal of the discriminator, when shown an instance from the true MNIST dataset, is to recognize those that are authentic.

Meanwhile, the generator is creating new, synthetic images that it passes to the discriminator. It does so in the hopes that they, too, will be deemed authentic, even though they are fake. The goal of the generator is to generate passable hand-written digits: to lie without being caught. The goal of the discriminator is to identify images coming from the generator as fake.
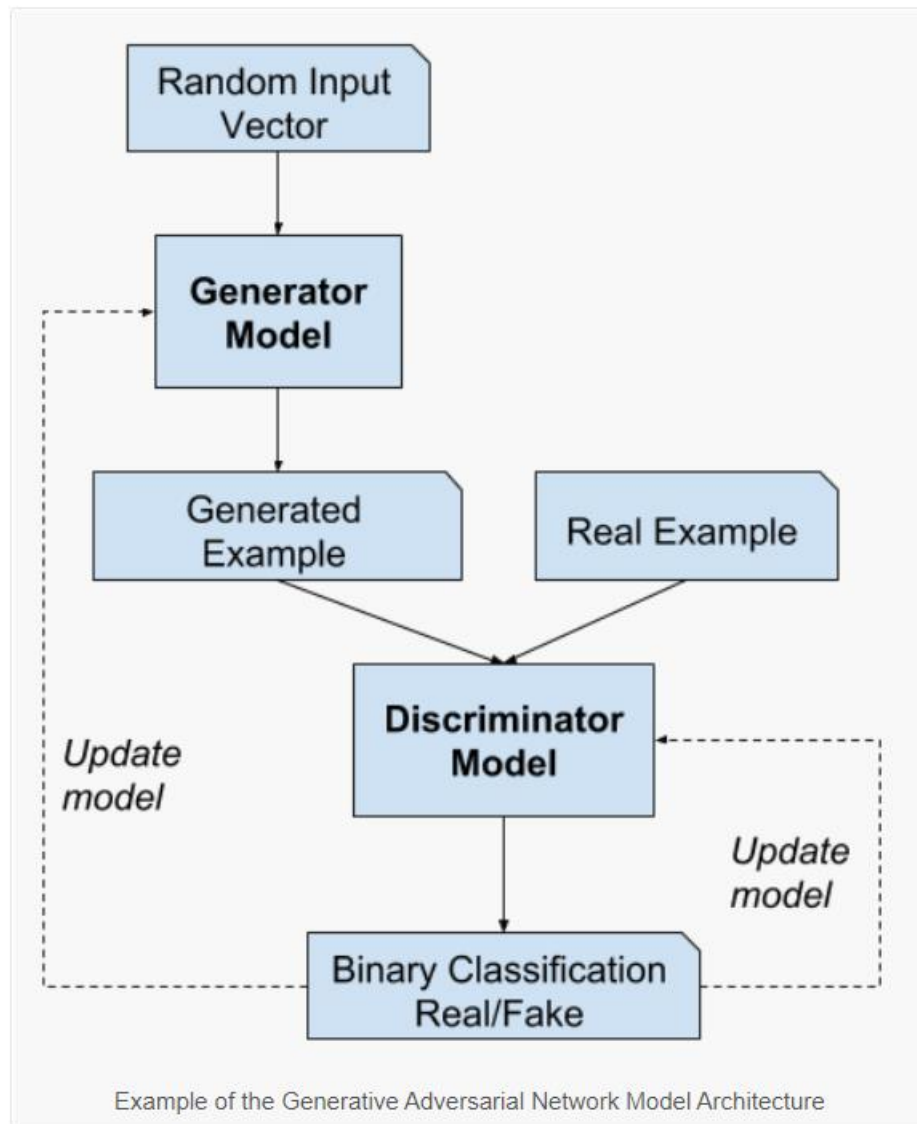
Here are the steps a GAN takes:

- The generator takes in random numbers and returns an image.

- This generated image is fed into the discriminator alongside a stream of images taken from the actual, ground-truth dataset.

- The discriminator takes in both real and fake images and returns probabilities, a number between 0 and 1, with 1 representing a prediction of authenticity and 0 representing fake.

So, you have a double feedback loop:

- The discriminator is in a feedback loop with the ground truth of the images, which we know.

- The generator is in a feedback loop with the discriminator.

You can think of a GAN as the opposition of a counterfeiter and a cop in a game of cat and mouse, where the counterfeiter is learning to pass false notes, and the cop is learning to detect them. Both are dynamic; i.e. the cop is in training, too (to extend the analogy, maybe the central bank is flagging bills that slipped through), and each side comes to learn the other's methods in a constant escalation.

Example of the Generative Adversarial Network Model Architecture

For MNIST, the discriminator network is a standard convolutional network that can categorize the images fed to it, a binomial classifier labeling images as real or fake. The generator is an inverse convolutional network, in a sense: While a standard convolutional classifier takes an image and downsamples it to produce a probability, the generator takes a vector of random noise and upsamples it to an image. The first throws away data through downsampling techniques like maxpooling, and the second generates new data.