

# Study Material - Youtube

## Document Information

- **Generated:** 2025-08-01 22:00:51
- **Source:** <https://youtu.be/eH9skKyqjJM>
- **Platform:** Youtube
- **Word Count:** 2,408 words
- **Estimated Reading Time:** ~12 minutes
- **Number of Chapters:** 8
- **Transcript Available:** Yes (analyzed from video content)

## Table of Contents

1. Bidirectional Generative Adversarial Networks (Bi-GANs) - Deep Understanding
  2. Practical Implementation in PyTorch
  3. At 09:04, the Generator class is defined.
  4. At 09:37, the Encoder class is defined.
  5. At 09:57, the Discriminator class is defined.
  6. No sigmoid! BCEWithLogitsLoss handles that
  7. Self-Assessment for This Video
  8. Key Takeaways from This Video
- 

## Video Overview

This video provides a comprehensive tutorial on **Bidirectional Generative Adversarial Networks (Bi-GANs)**. The lecture begins by establishing the foundational concepts of a standard “vanilla” GAN, highlighting its limitation in mapping from the image space back to the latent space. It then introduces the Bi-GAN architecture, which adds an **Encoder** network to create a two-way mapping between the latent space and the image space. The core of the tutorial involves a detailed mathematical and conceptual breakdown of the Bi-GAN components—the Generator, the new Encoder, and a modified Discriminator that evaluates joint distributions. The lecture culminates in a practical, step-by-step implementation of a Bi-GAN using PyTorch on the MNIST dataset, covering model definitions, the unique loss function, and the complete training loop.

## Learning Objectives

Upon completing this lecture, students will be able to: - **Contrast** the architecture of a standard GAN with that of a Bi-GAN. - **Understand the role of the Encoder** in a Bi-GAN for mapping images to the latent space. - **Explain how the Discriminator’s function is modified** in a Bi-GAN to assess joint (image, latent vector) pairs. - **Formulate the mathematical mappings** for the Generator ( $G_\theta$ ), Encoder ( $E_\phi$ ), and Discriminator ( $D_\omega$ ). - **Derive and interpret the Bi-GAN loss function** and the associated minimax game. - **Implement a Bi-GAN from scratch in PyTorch**, including defining the three network architectures and the corresponding training logic.

## Prerequisites

To fully benefit from this tutorial, students should have a solid understanding of: - **Fundamentals of Generative Adversarial Networks (GANs)**, including the adversarial training process. - **Core concepts of Deep Learning and Neural Networks**, such as linear layers, activation functions (ReLU, Sigmoid, Tanh), and backpropagation. - **Basic proficiency in Python and the PyTorch library**, including defining `nn.Module` classes, using optimizers, and creating data loaders. - **Foundational knowledge of probability and distributions**, particularly the Normal distribution.

## Key Concepts Covered

- Vanilla GAN vs. Bidirectional GAN (Bi-GAN)
  - Generator, Encoder, and Discriminator Networks
  - Latent Space ( $Z$ ) and Image Space ( $X$ )
  - Bidirectional Mapping ( $Z \leftrightarrow X$ )
  - Joint Distribution Discrimination
  - Bi-GAN Loss Function
  - PyTorch Implementation and Training Loop
- 

# Bidirectional Generative Adversarial Networks (Bi-GANs) - Deep Understanding

## Conceptual Framework: From One-Way to Two-Way Generation

### Intuitive Foundation: The Limitation of Standard GANs

A standard Generative Adversarial Network (GAN), often called a “vanilla” GAN, is exceptionally good at one task: learning a mapping from a simple, low-dimensional latent space ( $Z$ ) to a complex, high-dimensional data space, like the space of images ( $X$ ).

As introduced at **00:47**, this process can be visualized as follows: 1. A **Generator** network ( $G$ ) takes a random noise vector  $z$  (sampled from a known distribution like a Gaussian) as input. 2. It transforms this noise vector into a synthetic image  $\hat{x}$  that resembles the real images from the training dataset.

This is a **unidirectional** flow:  $Z \rightarrow X$ .

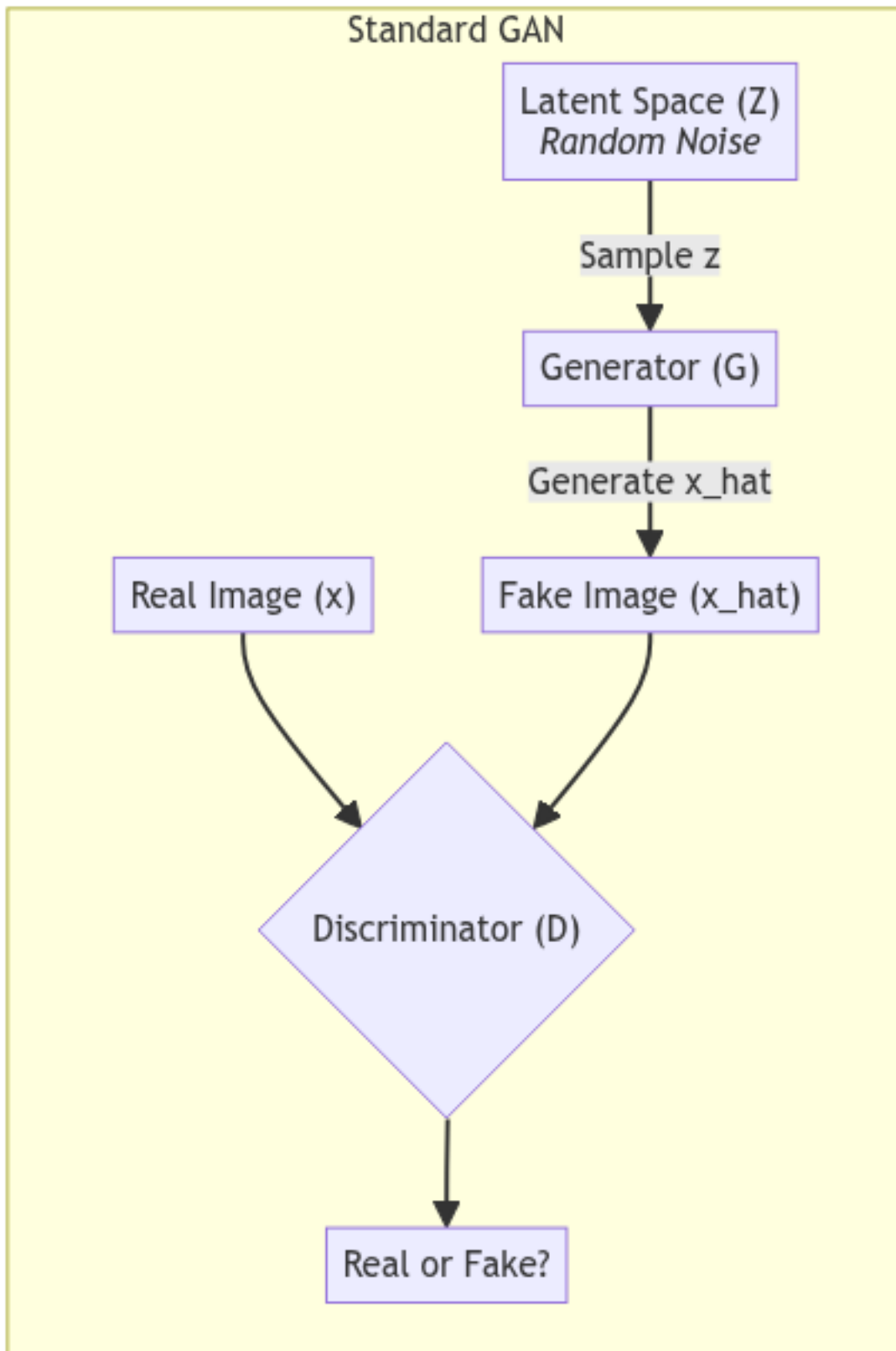


Figure 1: A simplified flowchart of a standard GAN, illustrating the one-way mapping from latent space  $Z$  to image space  $X$ .

However, this architecture lacks a crucial capability: **inversion** or **embedding**. Given a real image  $x$ , a standard GAN cannot tell us which latent vector  $z$  would produce it. There is no direct path from  $X$  back to  $Z$ . This is the problem that Bidirectional GANs (Bi-GANs) are designed to solve (01:38).

### The Bi-GAN Solution: Introducing the Encoder

A Bi-GAN introduces a third component, an **Encoder** network ( $E$ ), to create a two-way street between the latent and image spaces.

1. **Generator ( $G$ )**: Same as before, maps from latent space to image space ( $G : Z \rightarrow X$ ).
2. **Encoder ( $E$ )**: The new component, learns the reverse mapping from image space to latent space ( $E : X \rightarrow Z$ ).
3. **Discriminator ( $D$ )**: This network is modified to become a “joint discriminator.” Instead of just classifying images as real or fake, it now classifies **pairs** of (image, latent vector) as jointly “real” or “fake.”

This creates a complete, bidirectional loop, as visualized below.

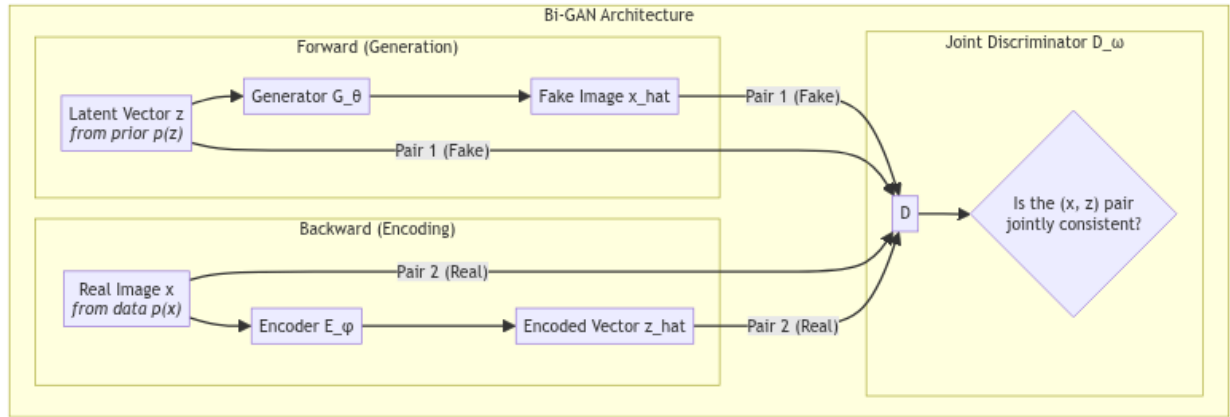


Figure 2: The architecture of a Bidirectional GAN (Bi-GAN). The Discriminator learns to distinguish between pairs generated by the Generator (Fake) and pairs generated by the Encoder (Real).

The core idea is that the discriminator forces the generator and encoder to learn inverse mappings of each other. If the generator creates a realistic image from a latent vector  $z$ , the encoder should be able to take that same image and map it back to the original  $z$ .

## Mathematical Analysis of Bi-GAN

### Network Mappings

The Bi-GAN framework consists of three distinct networks, each with a specific mapping function (01:58):

1. **Generator ( $G_\theta$ )**: Maps a latent vector  $z$  from the prior distribution  $p_z(z)$  to the image space  $X$ .

$$G_\theta : Z \rightarrow X$$

The generated image is  $\hat{x} = G_\theta(z)$ .

2. **Encoder ( $E_\phi$ )**: Maps a real image  $x$  from the data distribution  $p_x(x)$  to the latent space  $Z$ .

$$E_\phi : X \rightarrow Z$$

The encoded latent vector is  $\hat{z} = E_\phi(x)$ .

3. **Discriminator** ( $D_\omega$ ): This is a joint discriminator that takes a pair  $(x, z)$  and outputs a single scalar value indicating whether the pair is from the real or fake joint distribution.

$$D_\omega : X \times Z \rightarrow [0, 1]$$

## The Bi-GAN Objective Function

The training objective is a minimax game, similar to a standard GAN, but adapted for the joint distributions. The discriminator  $D$  tries to maximize the following objective function, while the generator  $G$  and encoder  $E$  try to minimize it (**05:48**).

The objective function  $L_{BiGAN}$  is:

$$\min_{G,E} \max_D L(G, E, D) = \mathbb{E}_{\mathbf{x} \sim p_x} [\log D_\omega(\mathbf{x}, E_\phi(\mathbf{x}))] + \mathbb{E}_{\mathbf{z} \sim p_z} [\log(1 - D_\omega(G_\theta(\mathbf{z}), \mathbf{z}))]$$

## Intuitive Breakdown of the Loss Function:

- **Term 1:**  $\mathbb{E}_{\mathbf{x} \sim p_x} [\log D_\omega(\mathbf{x}, E_\phi(\mathbf{x}))]$ 
  - This term deals with the **“real” pairs**.
  - We take a real image  $\mathbf{x}$  from the dataset.
  - We use the **Encoder**  $E_\phi$  to get its latent representation,  $E_\phi(\mathbf{x})$ .
  - The pair  $(\mathbf{x}, E_\phi(\mathbf{x}))$  is fed to the discriminator.
  - The discriminator  $D_\omega$  is trained to maximize this term, meaning it should output a value close to 1 (high probability of being a real pair). The encoder  $E_\phi$  is trained to help the discriminator by producing consistent latent codes.
- **Term 2:**  $\mathbb{E}_{\mathbf{z} \sim p_z} [\log(1 - D_\omega(G_\theta(\mathbf{z}), \mathbf{z}))]$ 
  - This term deals with the **“fake” pairs**.
  - We take a random latent vector  $\mathbf{z}$  from the prior distribution.
  - We use the **Generator**  $G_\theta$  to create a fake image,  $G_\theta(\mathbf{z})$ .
  - The pair  $(G_\theta(\mathbf{z}), \mathbf{z})$  is fed to the discriminator.
  - The discriminator  $D_\omega$  is trained to make  $D_\omega(G_\theta(\mathbf{z}), \mathbf{z})$  close to 0, which maximizes  $\log(1 - D_\omega(\dots))$ .
  - The generator  $G_\theta$  is trained to do the opposite: it wants to fool the discriminator by making  $D_\omega(G_\theta(\mathbf{z}), \mathbf{z})$  close to 1, which *minimizes* this term.

This adversarial process forces the distribution of encoded real images,  $p_\phi(z|x)$ , to match the prior distribution  $p_z(z)$ , and the distribution of generated images,  $p_\theta(x|z)$ , to match the real data distribution  $p_x(x)$ .

## Practical Implementation in PyTorch

The video provides a detailed code walkthrough for implementing a Bi-GAN on the MNIST dataset. The code is structured into several key components.

### Visual Elements from the Video

At **08:37**, the instructor switches to a Google Colab notebook to demonstrate the implementation. The notebook is titled `IITM_DGM_Bi-GAN.ipynb`.

### 1. Configuration and Data Loading (08:48)

Standard setup for a PyTorch project.

- **Libraries:** `torch`, `torch.nn`, `torchvision`, etc. are imported.
- **Hyperparameters:**
  - `device`: Set to “cuda” if available, otherwise “cpu”.
  - `batch_size`: 128

- latent\_dim: 50 (The dimensionality of the latent space  $Z$ ).
- epochs: 30
- lr: 2e-4 (Learning rate for the optimizers).
- **Dataset:** MNIST is loaded using `torchvision.datasets.MNIST`.
  - A transform is applied to convert images to tensors and flatten them into 784-dimensional vectors ( $28 \times 28 = 784$ ).

## 2. Model Architectures

Three separate neural networks are defined.

**Generator (09:04)** The Generator maps the 50-dimensional latent space to the 784-dimensional image space.

*# At 09:04, the Generator class is defined.*

```
class Generator(nn.Module):
    def __init__(self):
        super().__init__()
        self.net = nn.Sequential(
            nn.Linear(latent_dim, 1024),
            nn.ReLU(),
            nn.Linear(1024, 1024),
            nn.BatchNorm1d(1024),
            nn.ReLU(),
            nn.Linear(1024, 784),
            nn.Tanh() # Corrected from Sigmoid at 14:20
        )
    def forward(self, z):
        return self.net(z)
```

*Note: The instructor initially uses `nn.Sigmoid` but later corrects it to `nn.Tanh` at 14:20. `Tanh` scales the output to  $[-1, 1]$ , which is a common practice for image generation.*

**Encoder (09:37)** The Encoder performs the reverse mapping, from the 784-dimensional image space to the 50-dimensional latent space.

*# At 09:37, the Encoder class is defined.*

```
class Encoder(nn.Module):
    def __init__(self):
        super().__init__()
        self.net = nn.Sequential(
            nn.Linear(784, 1024),
            nn.ReLU(),
            nn.Linear(1024, 1024),
            nn.BatchNorm1d(1024),
            nn.ReLU(),
            nn.Linear(1024, latent_dim)
        )
    def forward(self, x):
        return self.net(x)
```

**Discriminator (09:57)** The Discriminator takes a concatenated pair of an image (784 dims) and a latent vector (50 dims) as input.

*# At 09:57, the Discriminator class is defined.*

```
class Discriminator(nn.Module):
```

```

def __init__(self):
    super().__init__()
    self.net = nn.Sequential(
        nn.Linear(784 + latent_dim, 1024),
        nn.ReLU(),
        nn.Linear(1024, 1024),
        nn.BatchNorm1d(1024),
        nn.ReLU(),
        nn.Linear(1024, 1)
        # No sigmoid! BCEWithLogitsLoss handles that
    )
def forward(self, x, z):
    xz = torch.cat((x, z), dim=1)
    return self.net(xz)

```

**Key Insight:** The input to the discriminator’s first linear layer is `784 + latent_dim`. The `forward` method explicitly concatenates the image `x` and latent vector `z` before passing them through the network.

### 3. Optimizers and Loss (10:33)

- **Loss Function:** `nn.BCEWithLogitsLoss()` is used, which combines a Sigmoid layer and the Binary Cross Entropy loss in one class for better numerical stability.
- **Optimizers:** Two Adam optimizers are created:
  1. `optimizer_D`: Manages the parameters of the Discriminator (`D.parameters()`).
  2. `optimizer_GE`: Manages the parameters of *both* the Generator and the Encoder (`list(G.parameters()) + list(E.parameters())`).

### 4. The Training Loop (10:57)

The training loop alternates between updating the discriminator and updating the generator/encoder.

#### Step 1: Train the Discriminator (11:12)

- The goal is to teach the discriminator to distinguish real pairs from fake pairs.
- **Real Pairs:** (`x_real`, `z_fake`) where `z_fake = E(x_real)`. The discriminator’s output `D_real` should be close to 1.
- **Fake Pairs:** (`x_fake`, `z_real`) where `x_fake = G(z_real)`. The discriminator’s output `D_fake` should be close to 0.
- The loss is calculated as `loss_D = bce_loss(D_real, real_labels) + bce_loss(D_fake, fake_labels)`.
- The gradients are computed, and the discriminator’s weights are updated.

#### Step 2: Train the Generator and Encoder (13:42)

- The goal is to fool the discriminator.
- The generator and encoder work together to make their respective pairs look “real” to the discriminator.
- **Generator’s Goal:** For the pair (`x_fake`, `z_real`), it wants the discriminator to output 1.
- **Encoder’s Goal:** For the pair (`x_real`, `z_fake`), it wants the discriminator to output 0.
- The loss is calculated by comparing the discriminator’s outputs to the *opposite* of the true labels: `loss_GE = bce_loss(D_real, fake_labels) + bce_loss(D_fake, real_labels)`.
- The gradients are computed, and the weights of *both* the generator and encoder are updated via `optimizer_GE`.

This process is repeated for a set number of epochs, gradually improving all three networks.

## Self-Assessment for This Video

1. **Question:** What is the primary architectural difference between a standard GAN and a Bi-GAN? Why is this change significant?
  - **Answer:** A Bi-GAN adds an Encoder network ( $E : X \rightarrow Z$ ) to the standard Generator-Discriminator setup. This is significant because it creates a bidirectional mapping, allowing the model not only to generate images from latent codes but also to find the latent code (embedding) for a given real image.
2. **Question:** Explain the inputs and output of the discriminator in a Bi-GAN. How does this differ from a standard GAN's discriminator?
  - **Answer:** In a standard GAN, the discriminator takes an image ( $x$ ) and outputs a probability. In a Bi-GAN, the discriminator takes a *pair* of an image and a latent vector ( $x, z$ ) and outputs a probability that this pair is jointly consistent (i.e., from the “real” joint distribution).
3. **Question:** Write down the simplified objective function for a Bi-GAN. Explain what the Generator/Encoder and the Discriminator are trying to achieve with respect to this function.
  - **Answer:** The objective is  $L = \mathbb{E}_{x \sim p_x} [\log D(x, E(x))] + \mathbb{E}_{z \sim p_z} [\log(1 - D(G(z), z))]$ . The Discriminator ( $D$ ) tries to maximize this value by correctly identifying real pairs (outputting 1) and fake pairs (outputting 0). The Generator ( $G$ ) and Encoder ( $E$ ) work together to minimize this value by fooling the discriminator.
4. **Application Exercise:** In the provided PyTorch code, if the `latent_dim` was changed to 100, what other parts of the code would need to be modified?
  - **Answer:**
    1. In the `Generator` class, the first `nn.Linear` layer's input dimension would change from `latent_dim` to 100.
    2. In the `Encoder` class, the last `nn.Linear` layer's output dimension would change to 100.
    3. In the `Discriminator` class, the first `nn.Linear` layer's input dimension would change from `784 + latent_dim` to `784 + 100`.

---

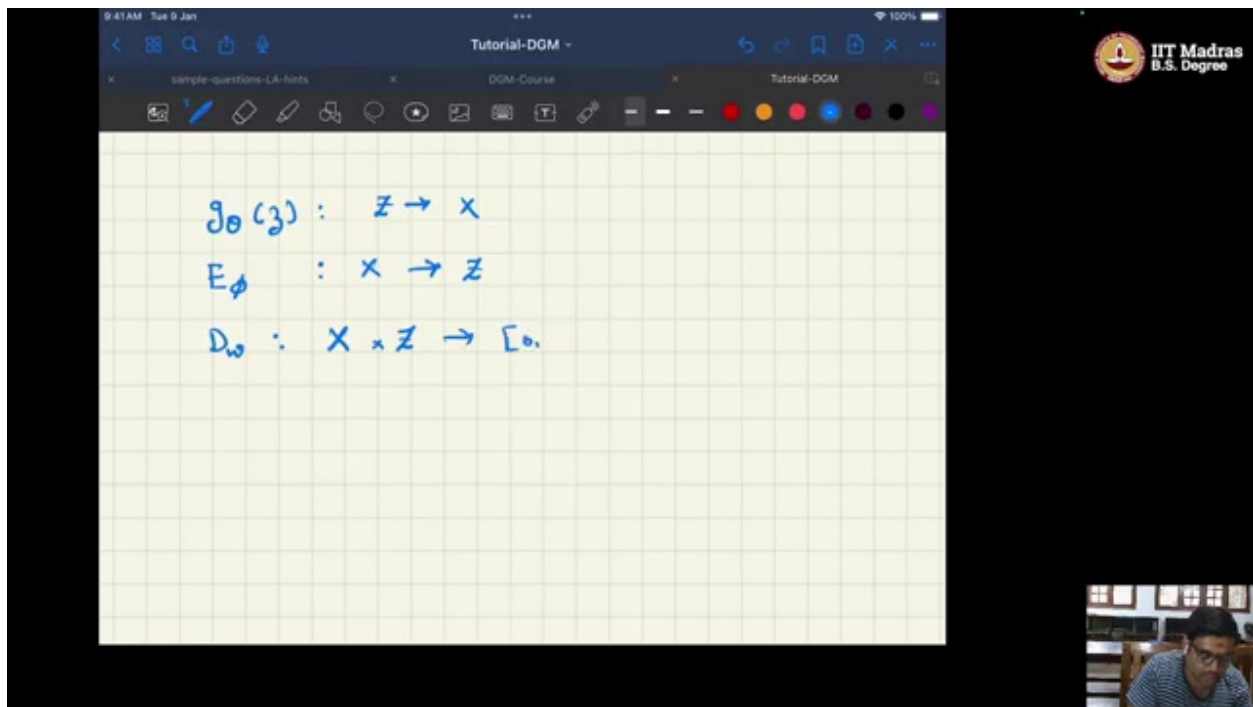
## Key Takeaways from This Video

- **Bi-GANs Enable Inversion:** They solve a key limitation of standard GANs by learning an encoder that maps data back to the latent space, which is useful for tasks like feature extraction and image retrieval.
- **The Discriminator's Role is Expanded:** The discriminator in a Bi-GAN is more powerful as it learns to model the joint probability distribution of images and their latent representations, not just the marginal distribution of images.
- **Training Involves a Three-Player Game:** The Generator, Encoder, and Discriminator are trained in an adversarial process. The Generator and Encoder team up to fool the Discriminator.
- **Implementation Requires Careful Pairing:** The key to implementing a Bi-GAN is correctly forming the (image, latent vector) pairs to feed into the discriminator for both the real and fake cases and using the appropriate labels for the loss calculation.

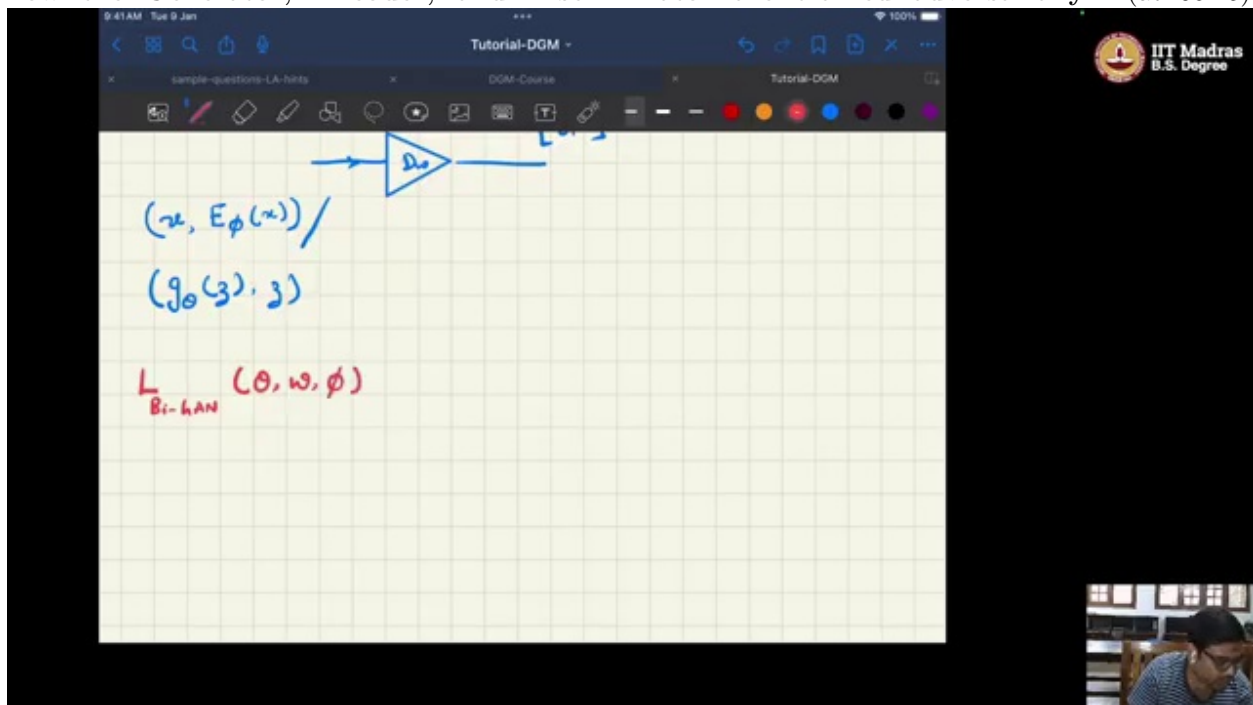
## Visual References

A diagram comparing the architecture of a standard GAN with a Bi-GAN. This visual is crucial for understanding the core innovation: the addition of an Encoder to create a two-way mapping between the latent space ( $Z$ ) and image space ( $X$ ). (at 02:30):





The slide presenting the mathematical formulation of the Bi-GAN loss function. This shows the complete minimax objective, which is essential for understanding how the Generator, Encoder, and Discriminator are trained adversarially. (at 06:15):



The PyTorch code implementation for the Discriminator class. This is a key practical example showing how the network is structured to take joint (image, latent vector) pairs as input, which is a fundamental difference from a standard GAN. (at 09:57):

The screenshot shows a Jupyter Notebook interface with the file name 'IITM\_DGM\_BI-GAN.ipynb'. The code defines a Discriminator class that takes a joint input of real data (x) and generated data (z). The class is initialized with a sequential network consisting of a linear layer (784 + latent\_dim to 1024), a ReLU activation, another linear layer (1024 to 1024), a Batch Normalization layer, a third ReLU activation, and a final linear layer (1024 to 1) for the output. The forward method concatenates the inputs x and z along the first dimension and passes them through the network.

```

11 nn.Linear(1024, latent_dim)
12 nn.ReLU(),
13 nn.Linear(1024, latent_dim)
14 )
15
16 def forward(self, x):
17     return self.net(x)
18
19
20 # Discriminator (Joint: x, z)
21 #
22 class Discriminator(nn.Module):
23     def __init__(self):
24         super().__init__()
25         self.net = nn.Sequential(
26             nn.Linear(784 + latent_dim, 1024),
27             nn.ReLU(),
28             nn.Linear(1024, 1024),
29             nn.BatchNorm1d(1024),
30             nn.ReLU(),
31             nn.Linear(1024, 1) # No sigmoid! BCEWithLogitsLoss handles that
32         )
33
34     def forward(self, x, z):
35         xz = torch.cat([x, z], dim=1)
36         return self.net(xz)

```

A final summary slide or concept map listing the key takeaways. This screenshot would consolidate the main concepts, such as the role of the Encoder, the modified Discriminator, and the bidirectional mapping, for effective review. (at 11:45):

The screenshot shows the training loop for the Discriminator. It iterates over epochs and batches of real data. For each batch, it samples a latent vector z from a prior distribution. It then trains the Discriminator by evaluating it on both real data (x\_real) and generated data (x\_fake). The loss is calculated as the sum of binary cross-entropy losses for both real and fake samples.

```

4 for epoch in range(1, epochs + 1):
5     for i, (x_real, _) in enumerate(train_loader):
6         x_real = x_real.to(device)
7         batch_size = x_real.size(0)
8
9         # Sample z from prior
10        z_real = torch.randn(batch_size, latent_dim, device=device)
11
12        #
13        # 1. Train Discriminator
14        #
15        for _ in range(p):
16            G.eval()
17            E.eval()
18            D.train()
19
20            x_fake = G(z_real).detach()
21            z_fake = E(x_real).detach()
22
23            D_real = D(x_real, z_fake)
24            D_fake = D(x_fake, z_real)
25
26            label_real = torch.ones_like(D_real)
27            label_fake = torch.zeros_like(D_fake)
28
29            loss_D = bce_loss(D_real, label_real) + bce_loss(D_fake, label_fake)
30
31            # Backpropagate D loss

```