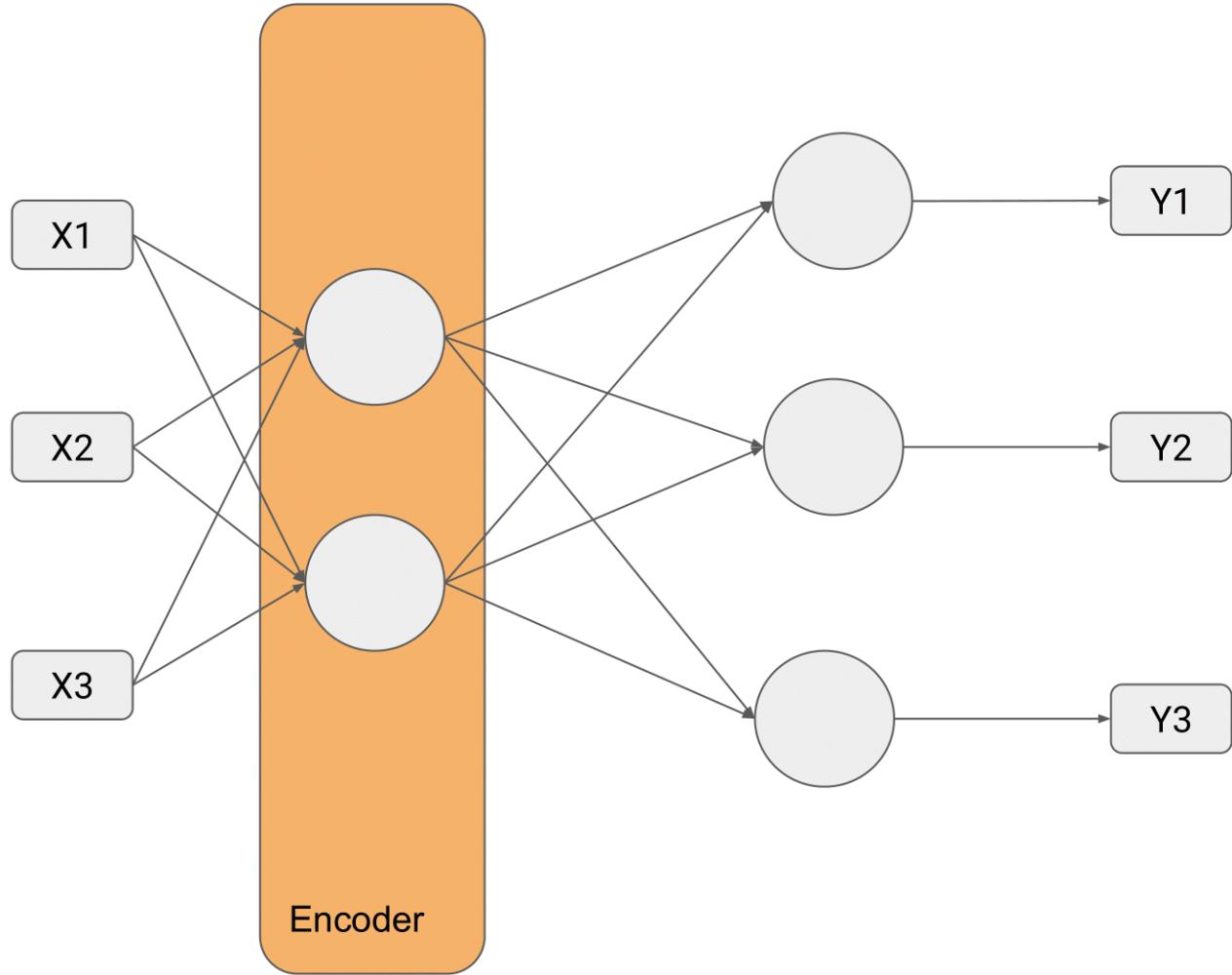
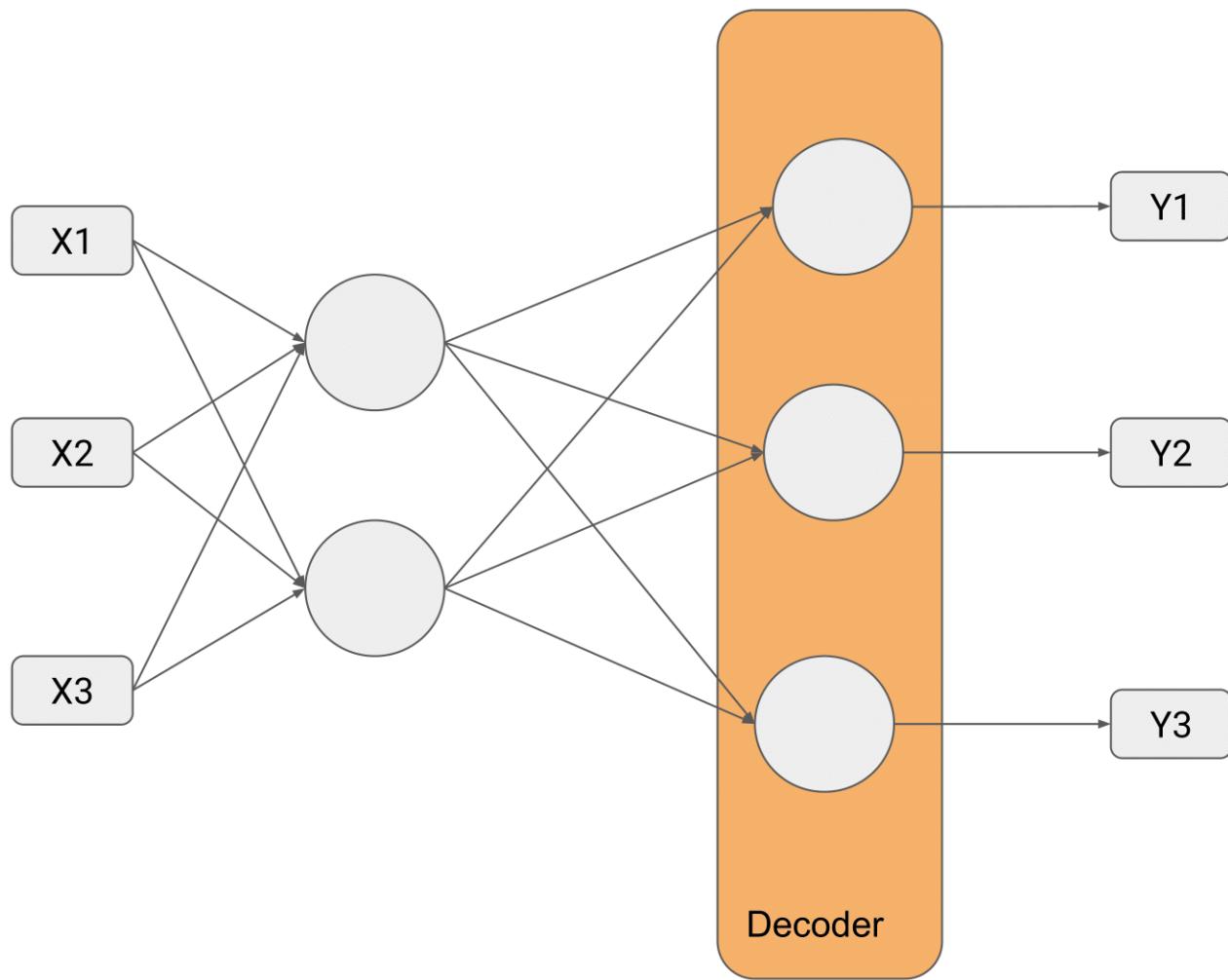
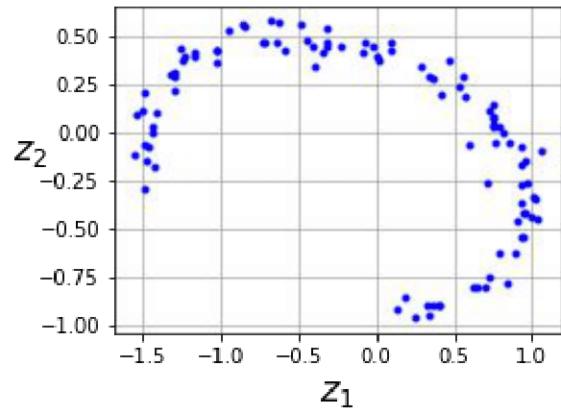
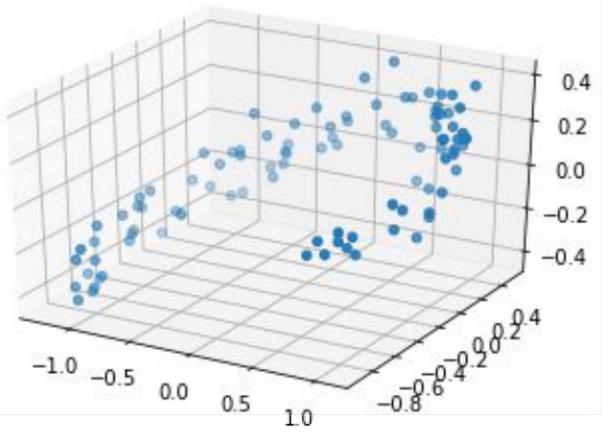


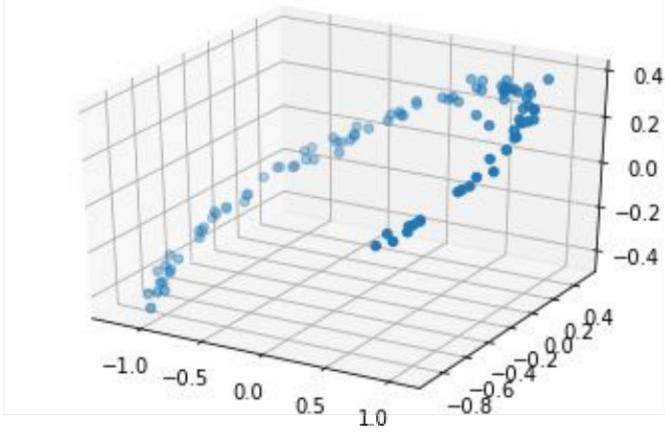
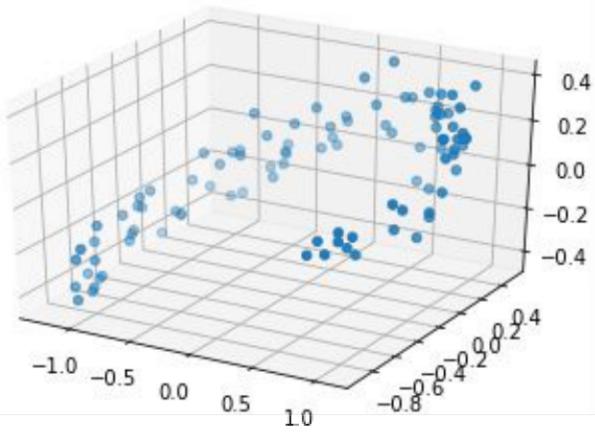
What are AutoEncoders?

- Neural networks capable of learning dense representations of input data without supervision
 - Training data is not labelled
- Useful for dimensionality reduction and for visualization
- Can be used to generate new data that resembles input data
- In practice they
 - Copy input to output
 - They learn efficient ways to represent data

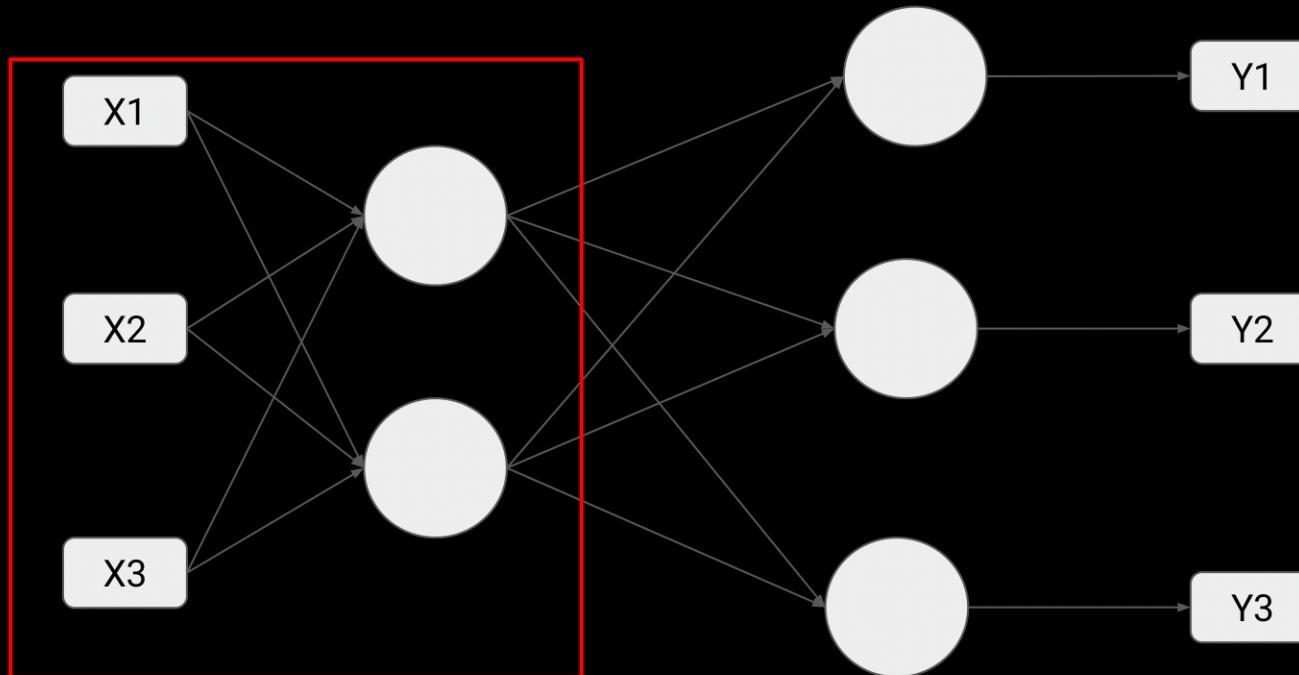








```
encoder = keras.models.Sequential([keras.layers.Dense(2, input_shape=[3])])  
  
decoder = keras.models.Sequential([keras.layers.Dense(3, input_shape=[2])])  
  
autoencoder = keras.models.Sequential([encoder, decoder])  
  
autoencoder.compile(loss="mse", optimizer=keras.optimizers.SGD(lr=1.5))
```

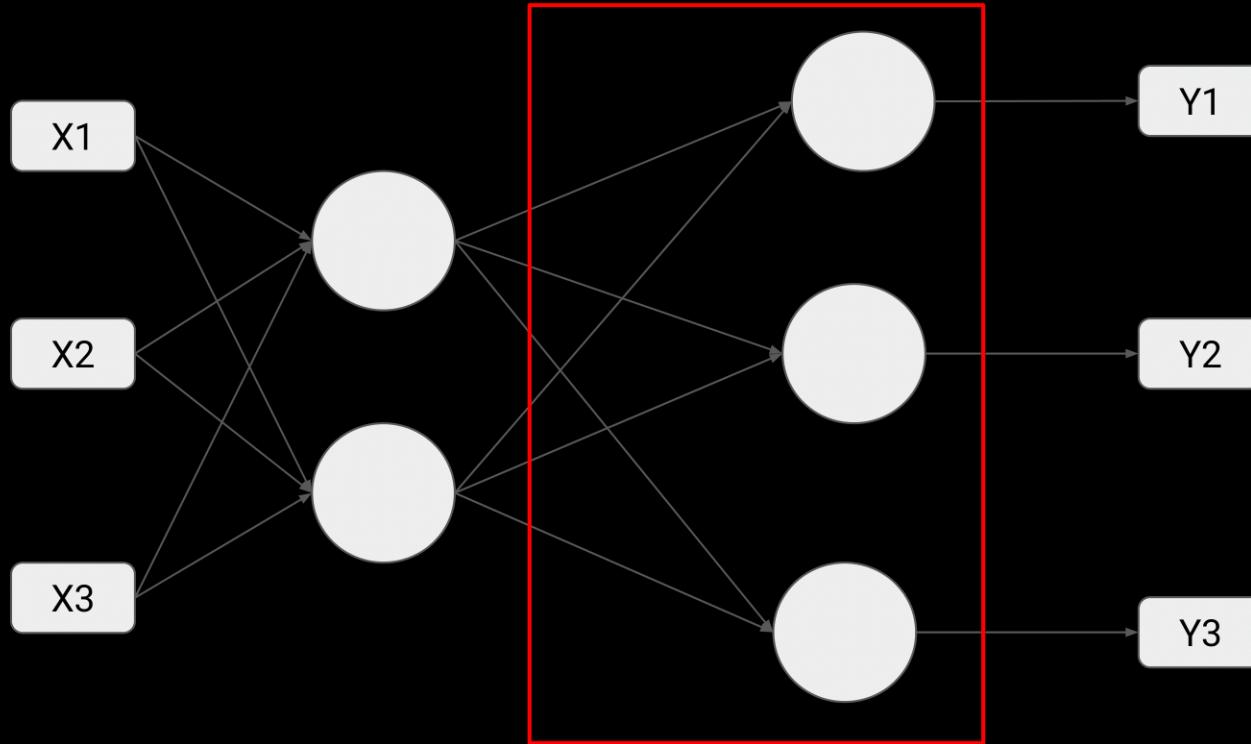


```
encoder = keras.models.Sequential([keras.layers.Dense(2, input_shape=[3])])
```

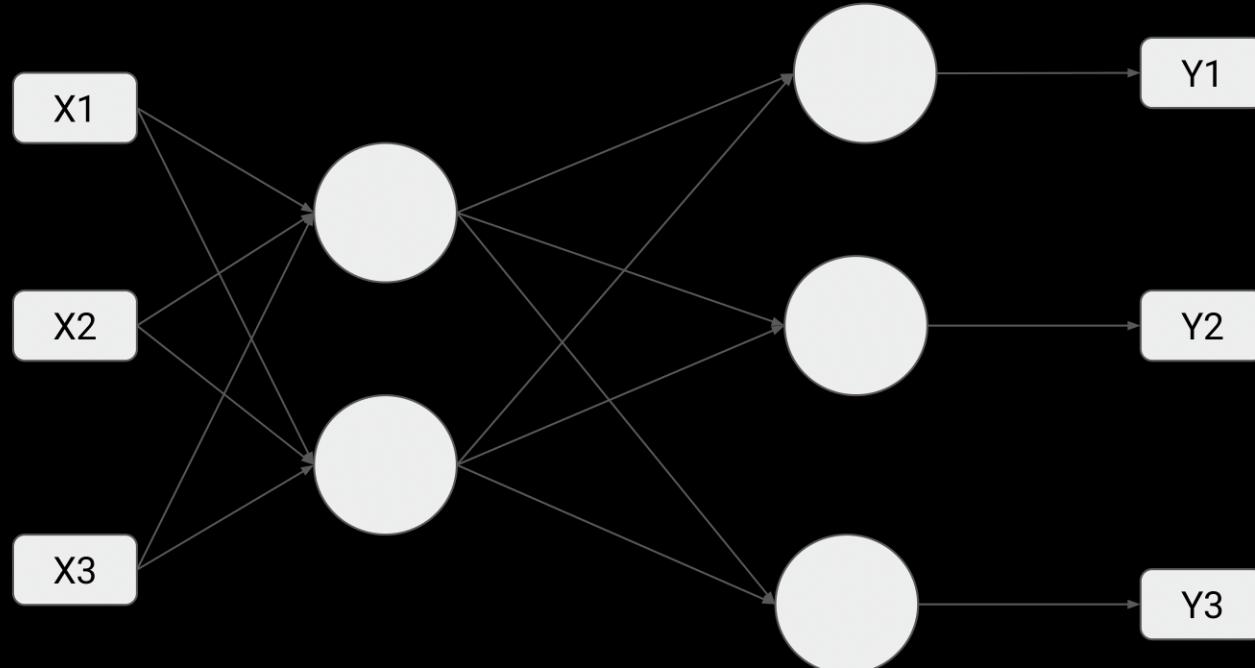
```
decoder = keras.models.Sequential([keras.layers.Dense(3, input_shape=[2])])
```

```
autoencoder = keras.models.Sequential([encoder, decoder])
```

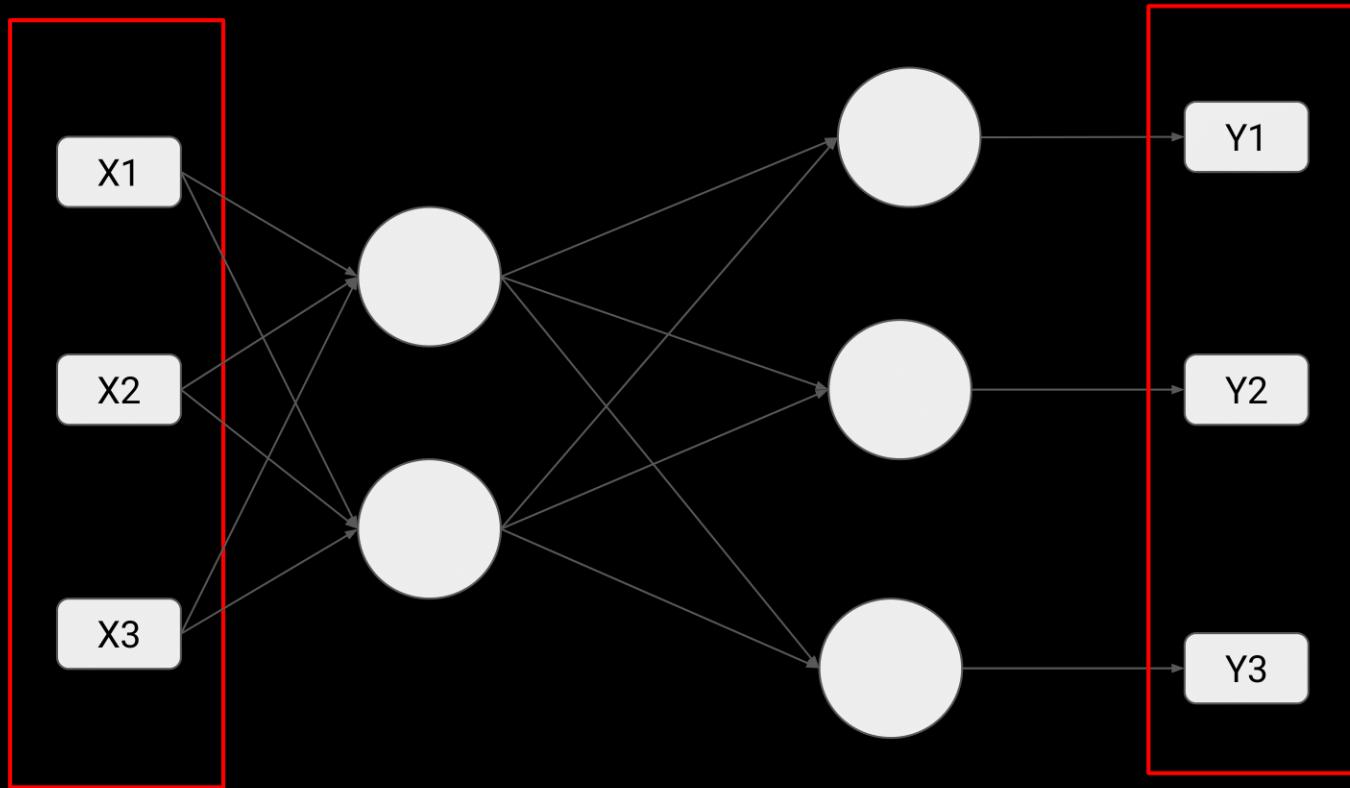
```
autoencoder.compile(loss="mse", optimizer=keras.optimizers.SGD(lr=1.5))
```



```
encoder = keras.models.Sequential([keras.layers.Dense(2, input_shape=[3])])  
decoder = keras.models.Sequential([keras.layers.Dense(3, input_shape=[2])])  
  
autoencoder = keras.models.Sequential([encoder, decoder])  
  
autoencoder.compile(loss="mse", optimizer=keras.optimizers.SGD(lr=1.5))
```



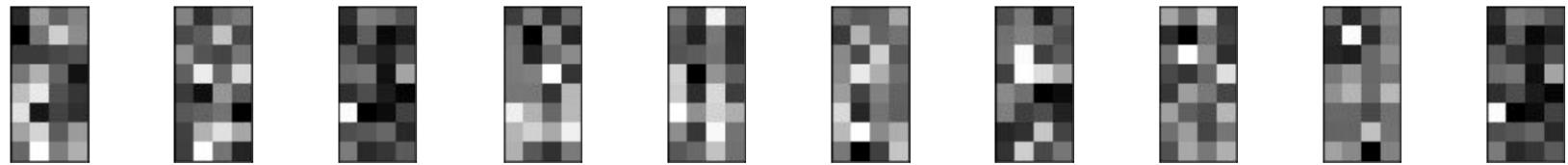
```
history = autoencoder.fit(X_train, X_train, epochs=200)
```



```
codings = encoder.predict(data)
```

```
decodings = decoder.predict(codings)
```

5 3 8 8 7 5 1 1 9 8



5 3 8 8 7 5 1 1 9 8

```
inputs = tf.keras.layers.Input(shape=(784,))

def simple_autoencoder():
    encoder = tf.keras.layers.Dense(units=32, activation='relu')(inputs)
    decoder = tf.keras.layers.Dense(units=784, activation='sigmoid')(encoder)
    return encoder, decoder

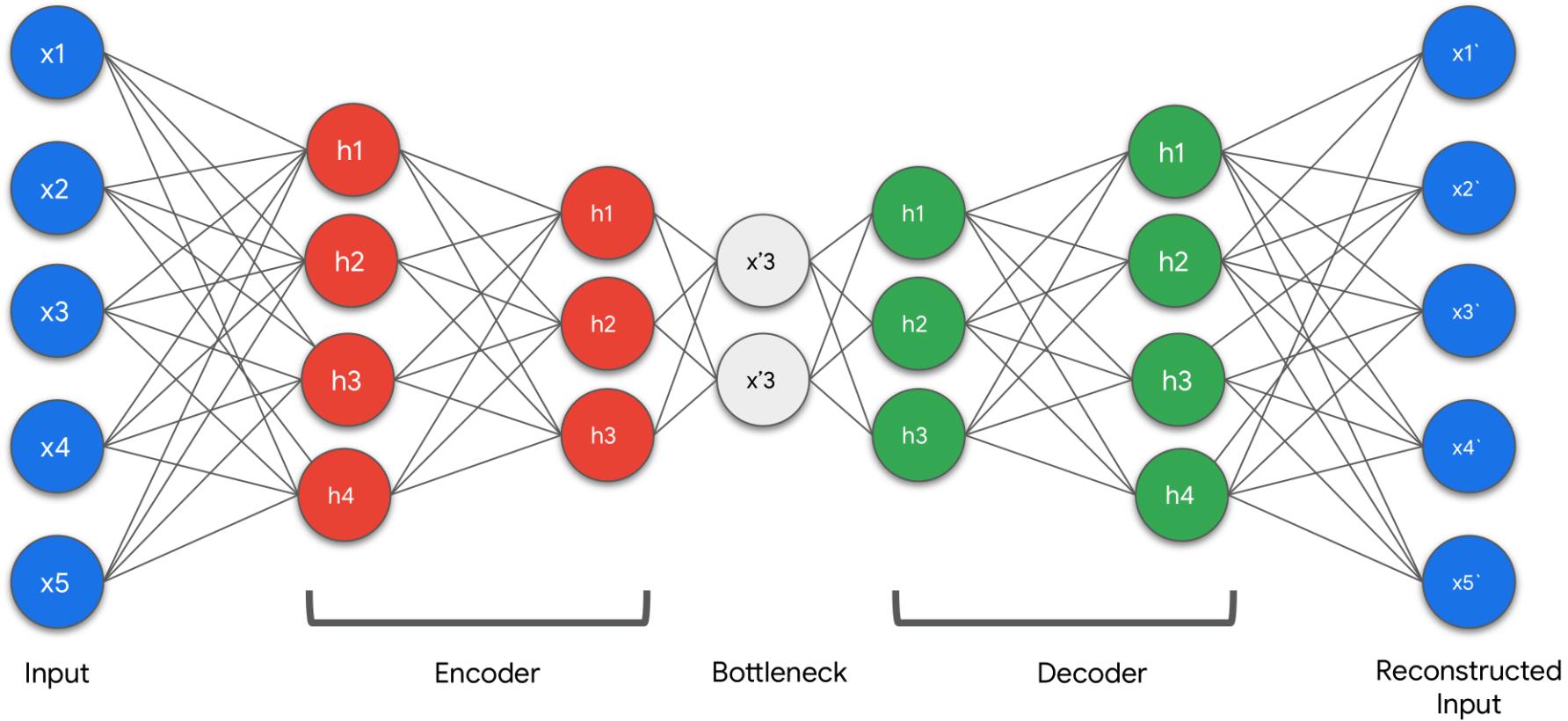
encoder_output, decoder_output = simple_autoencoder()

encoder_model = tf.keras.Model(inputs=inputs, outputs=encoder_output)

autoencoder_model = tf.keras.Model(inputs=inputs, outputs=decoder_output)
```

```
autoencoder_model.compile(  
    optimizer=tf.keras.optimizers.Adam(),  
    loss='binary_crossentropy')
```

Stacked Auto-Encoders



```
inputs = tf.keras.layers.Input(shape=(784,))

def deep_autoencoder():
    encoder = tf.keras.layers.Dense(units=128, activation='relu')(inputs)
    encoder = tf.keras.layers.Dense(units=64, activation='relu')(encoder)
    encoder = tf.keras.layers.Dense(units=32, activation='relu')(encoder)

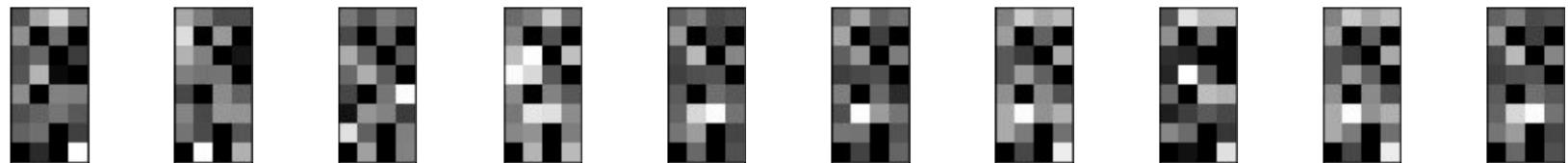
    decoder = tf.keras.layers.Dense(units=64, activation='relu')(encoder)
    decoder = tf.keras.layers.Dense(units=128, activation='relu')(decoder)
    decoder = tf.keras.layers.Dense(units=784, activation='sigmoid')(decoder)

    return encoder, decoder

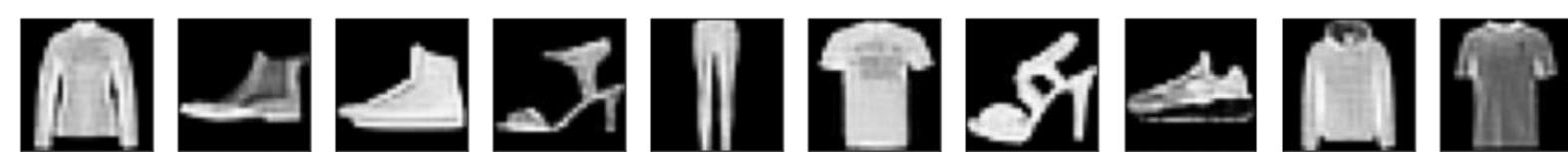
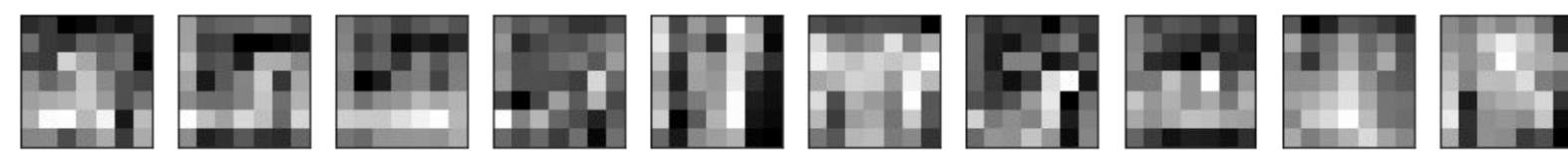
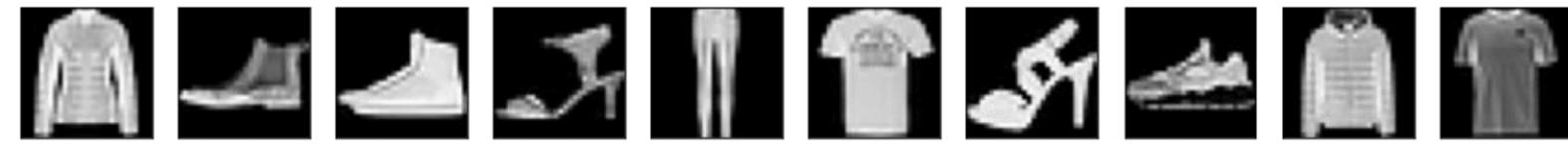
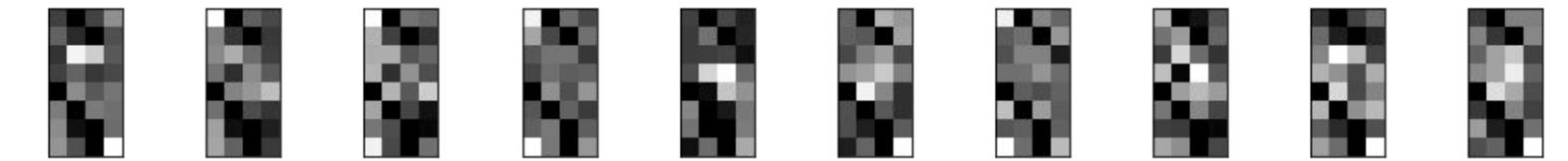
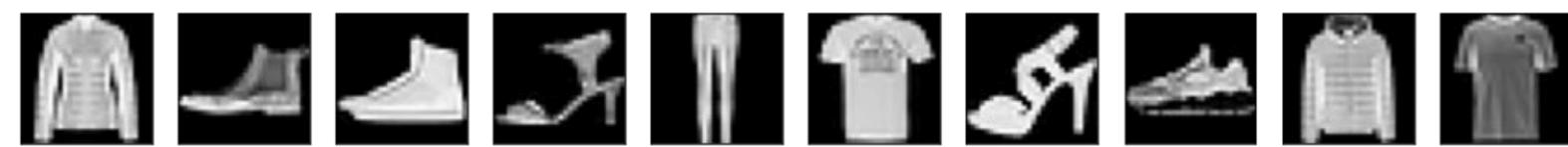
deep_encoder_output, deep_autoencoder_output = deep_autoencoder()

deep_encoder_model = tf.keras.Model(inputs=inputs, outputs=deep_encoder_output)
deep_autoencoder_model = tf.keras.Model(inputs=inputs, outputs=deep_autoencoder_output)
```

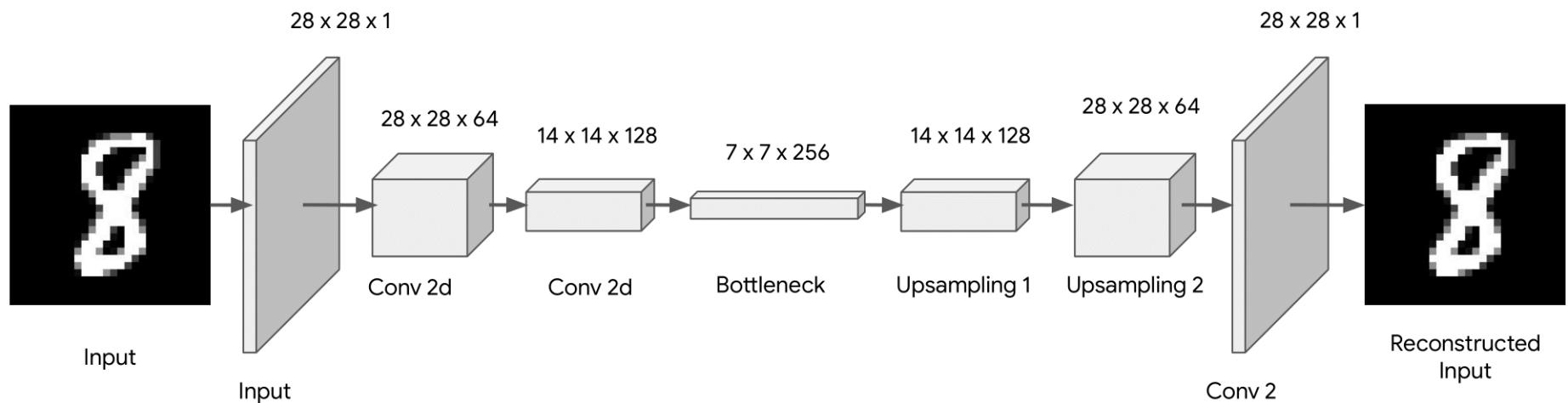
1 4 5 8 6 6 6 1 6 6



1 4 5 8 6 6 6 1 6 6



Convolutional Auto-Encoders



```
def encoder(inputs):

    conv_1 = tf.keras.layers.Conv2D(filters=64, kernel_size=(3,3),
                                    activation='relu', padding='same')(inputs)

    max_pool_1 = tf.keras.layers.MaxPooling2D(pool_size=(2,2))(conv_1)

    conv_2 = tf.keras.layers.Conv2D(filters=128, kernel_size=(3,3),
                                    activation='relu', padding='same')(max_pool_1)

    max_pool_2 = tf.keras.layers.MaxPooling2D(pool_size=(2,2))(conv_2)

    return max_pool_2
```

```
def bottle_neck(inputs):
    bottle_neck = tf.keras.layers.Conv2D(filters=256, kernel_size=(3,3),
                                         activation='relu', padding='same')(inputs)

    encoder_visualization = tf.keras.layers.Conv2D(filters=1, kernel_size=(3,3),
                                                   activation='sigmoid',
                                                   padding='same')(bottle_neck)

    return bottle_neck, encoder_visualization
```

```
def decoder(inputs):
    conv_1 = tf.keras.layers.Conv2D(filters=128, kernel_size=(3,3),
                                  activation='relu', padding='same')(inputs)
    up_sample_1 = tf.keras.layers.UpSampling2D(size=(2,2))(conv_1)

    conv_2 = tf.keras.layers.Conv2D(filters=64, kernel_size=(3,3),
                                  activation='relu', padding='same')(up_sample_1)
    up_sample_2 = tf.keras.layers.UpSampling2D(size=(2,2))(conv_2)

    conv_3 = tf.keras.layers.Conv2D(filters=1, kernel_size=(3,3),
                                  activation='sigmoid',
                                  padding='same')(up_sample_2)

    return conv_3
```

```
def convolutional_auto_encoder():
    inputs = tf.keras.layers.Input(shape=(28, 28, 1,))

    encoder_output = encoder(inputs)

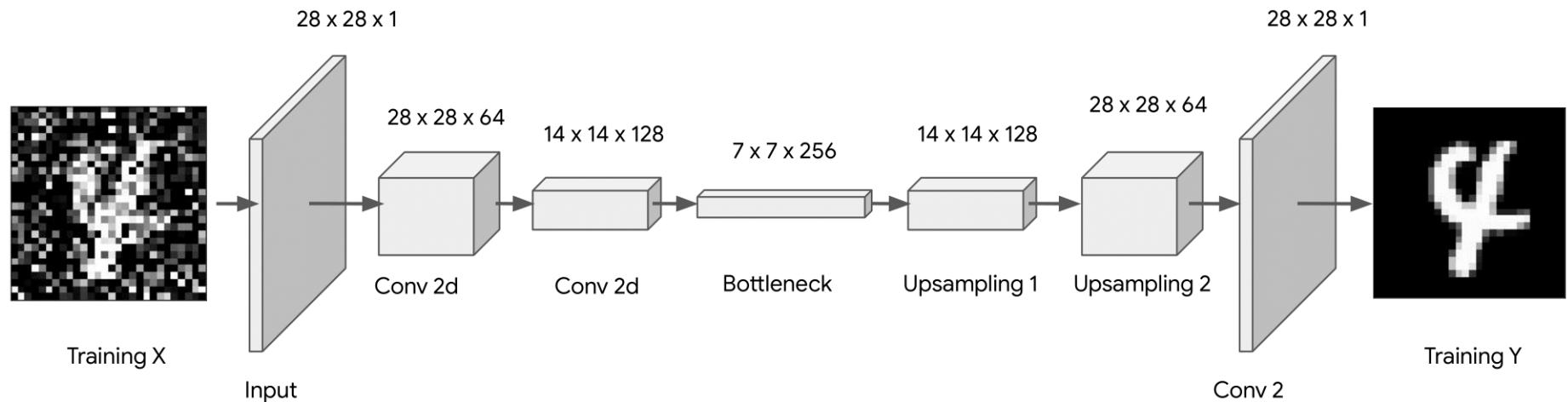
    bottleneck_output, encoder_visualization = bottle_neck(encoder_output)

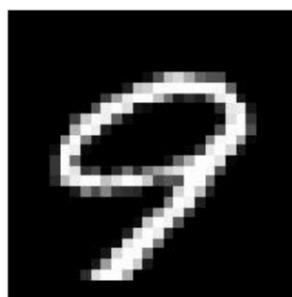
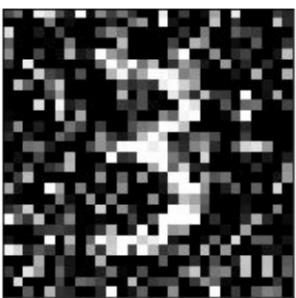
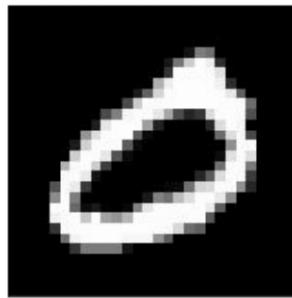
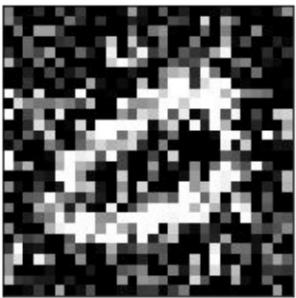
    decoder_output = decoder(bottleneck_output)

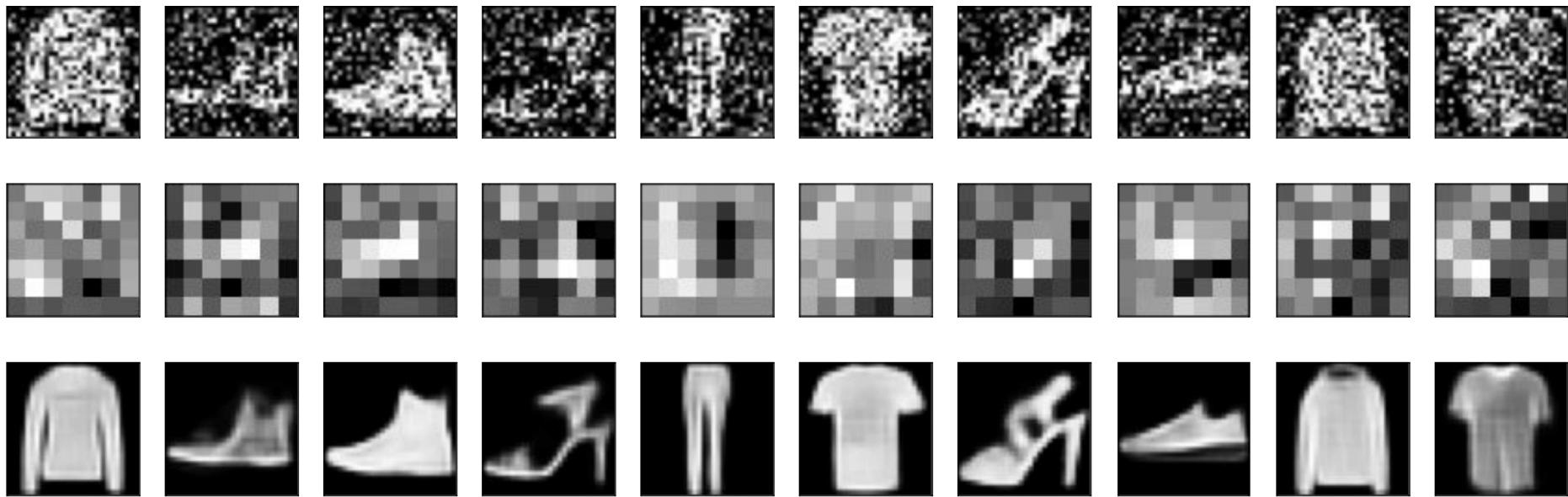
    model = tf.keras.Model(inputs =inputs, outputs=decoder_output)
    encoder_model = tf.keras.Model(inputs=inputs, outputs=encoder_visualization)

    return model, encoder_model
```

Convolutional Auto-Encoders



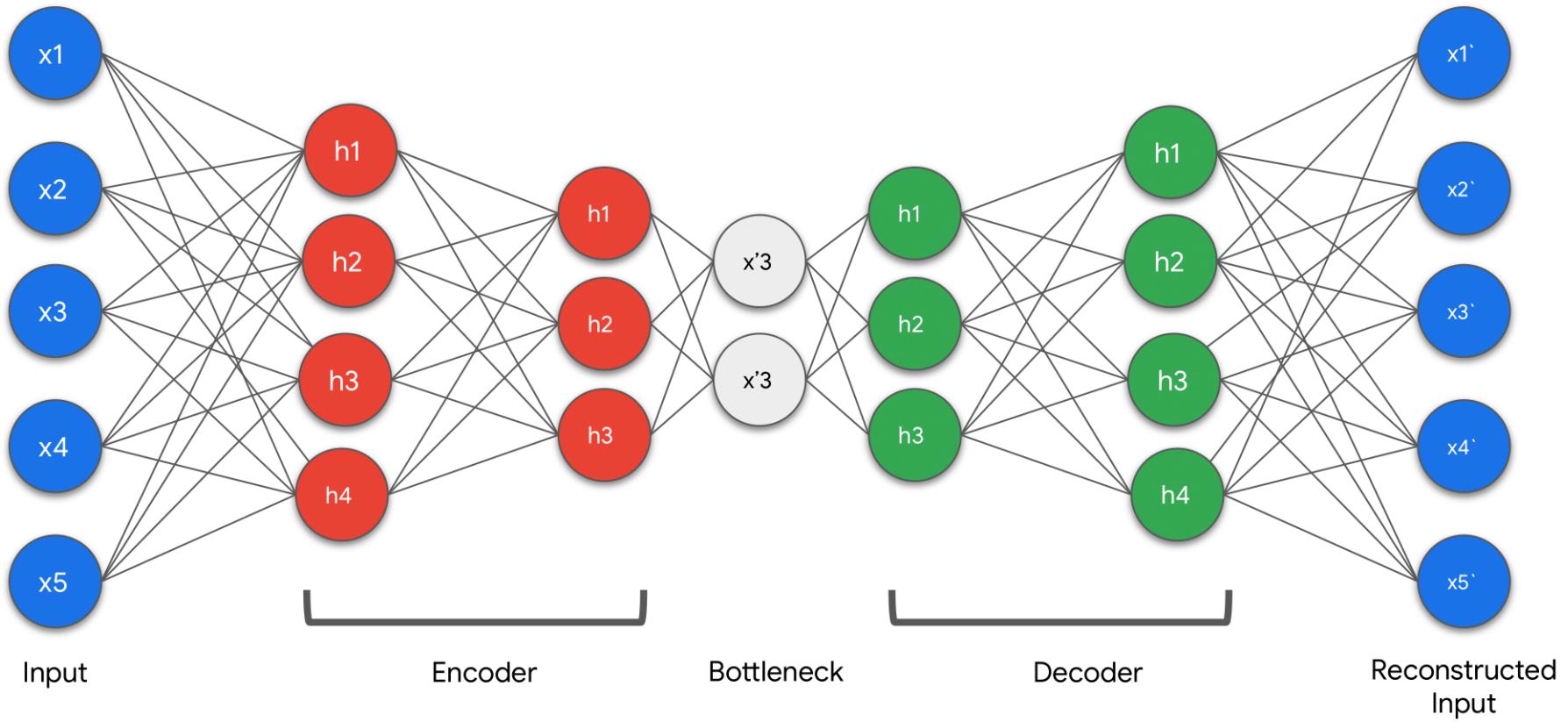


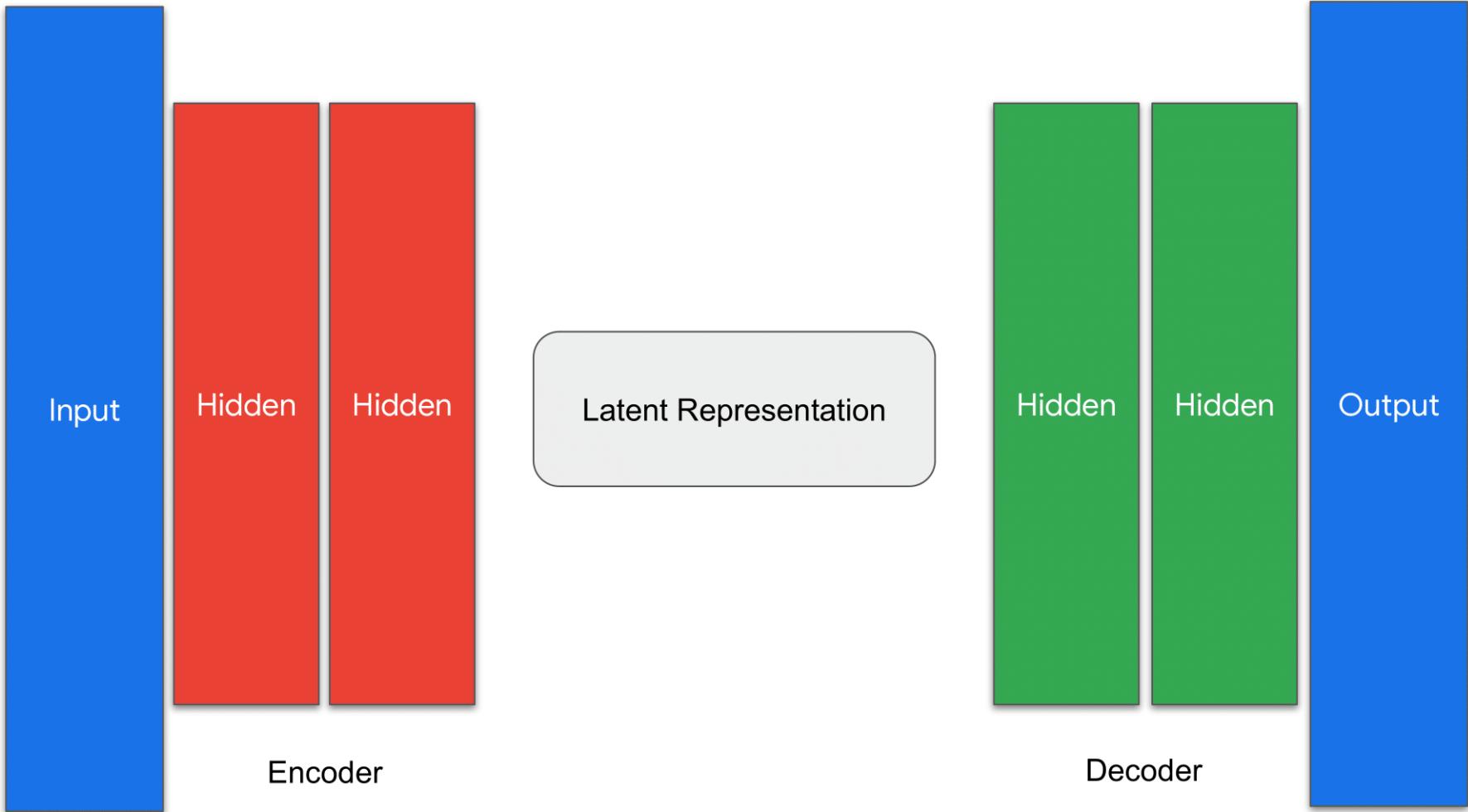


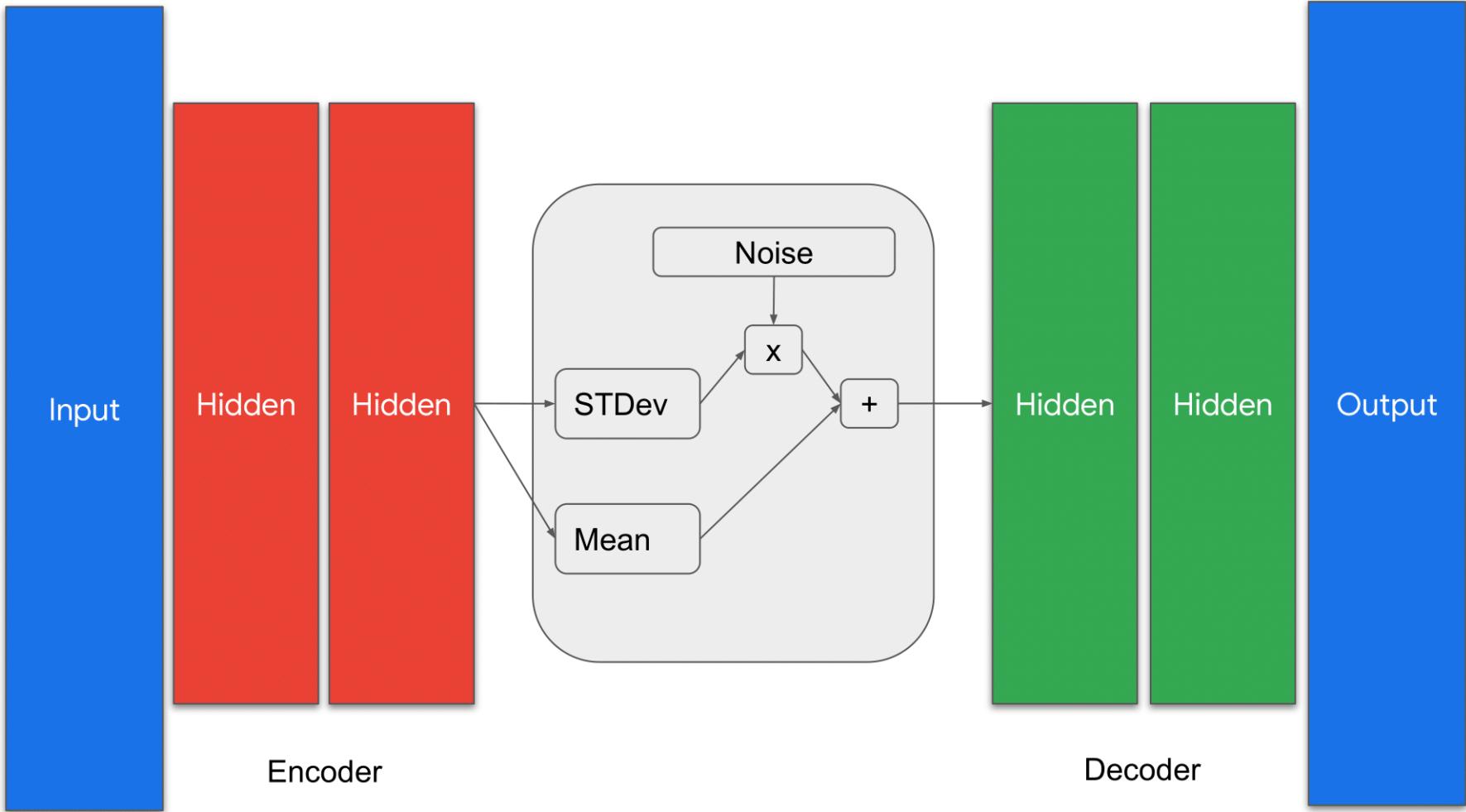
```
def map_image_with_noise(image, label):
    noise_factor = 0.5
    image = tf.cast(image, dtype=tf.float32)
    image = image / 255.0

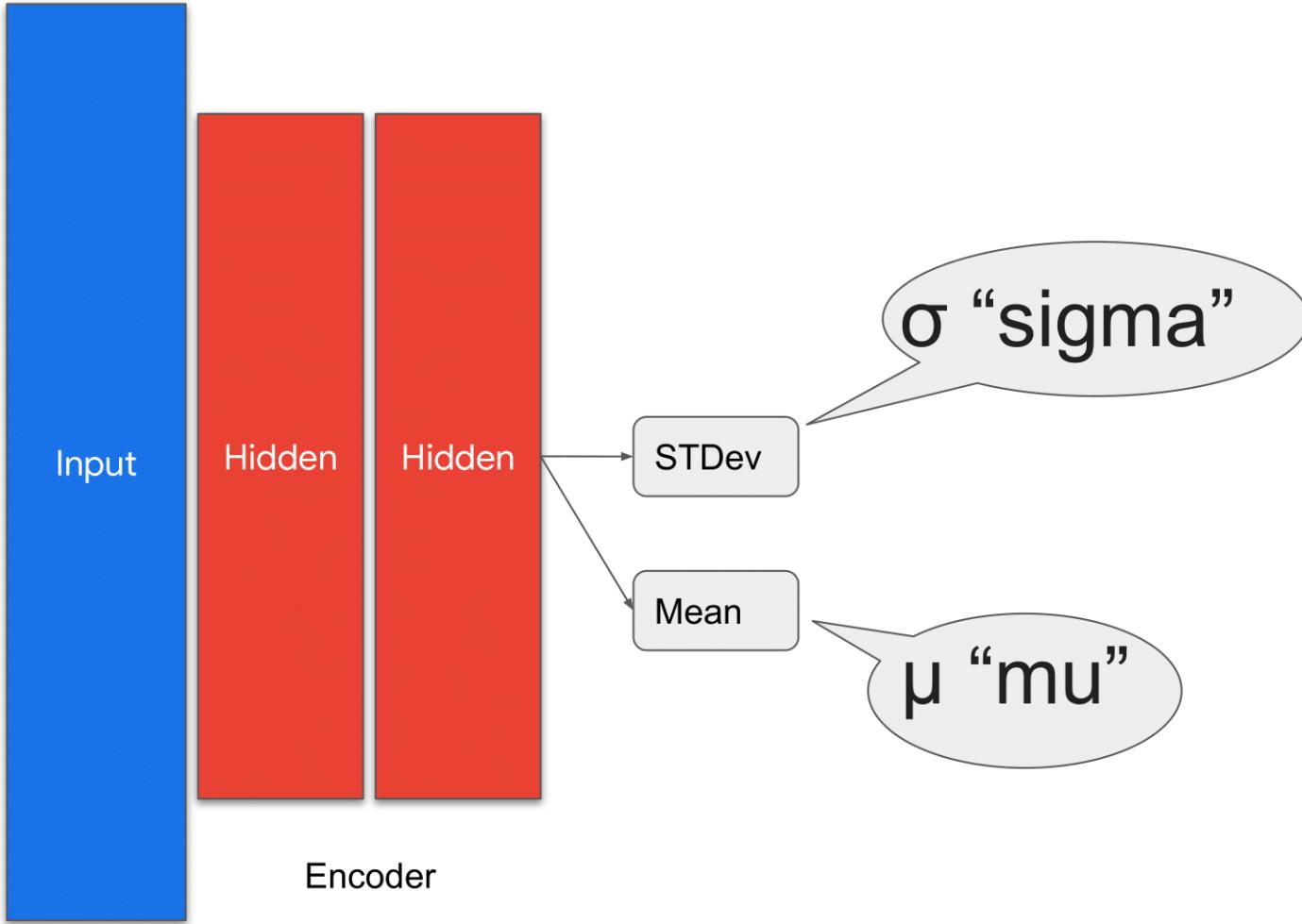
    factor = noise_factor * tf.random.normal(shape=image.shape)
    image_noisy = image + factor
    image_noisy = tf.clip_by_value(image_noisy, 0.0, 1.0)

    return image_noisy, image
```









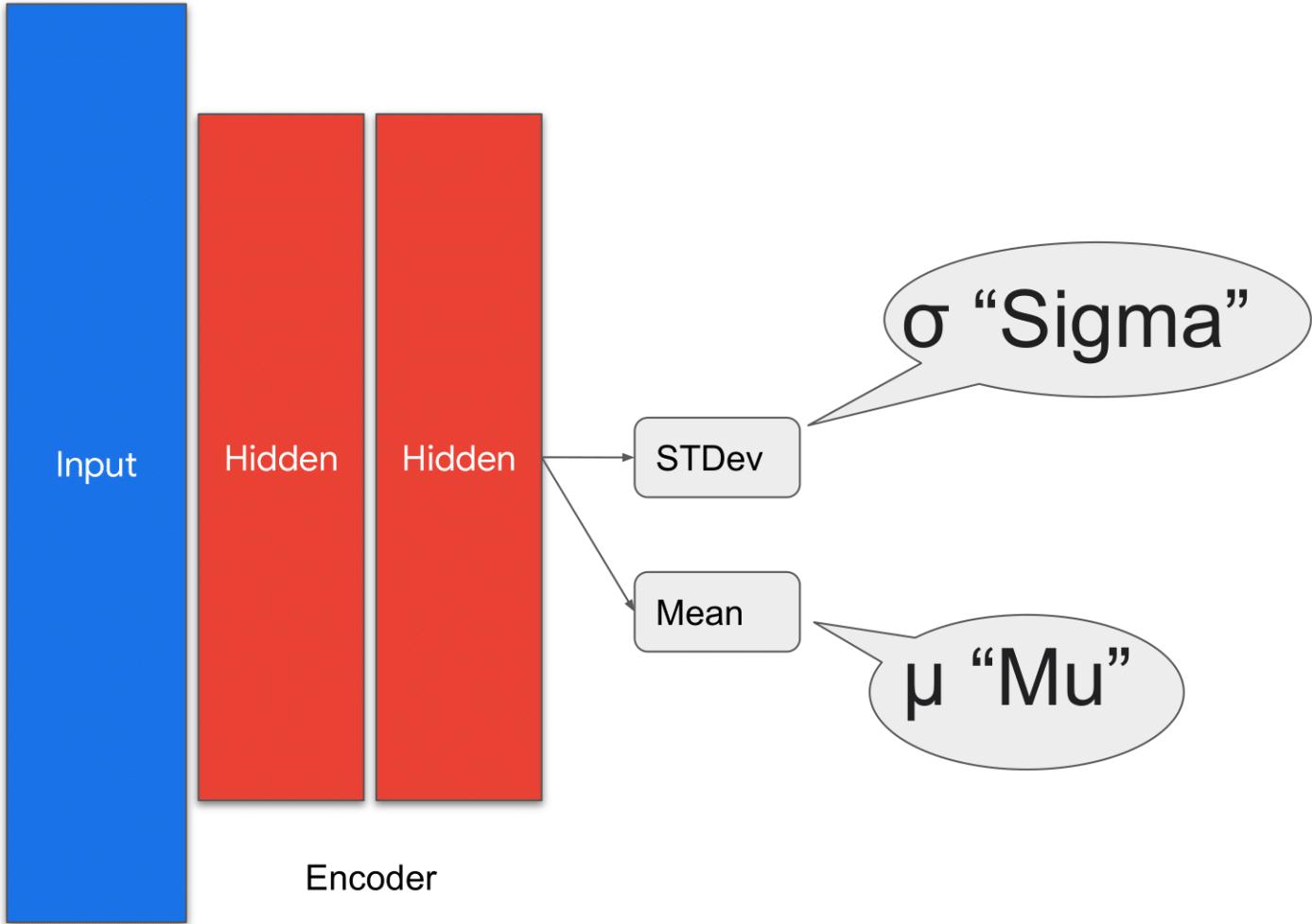
Probability Distribution

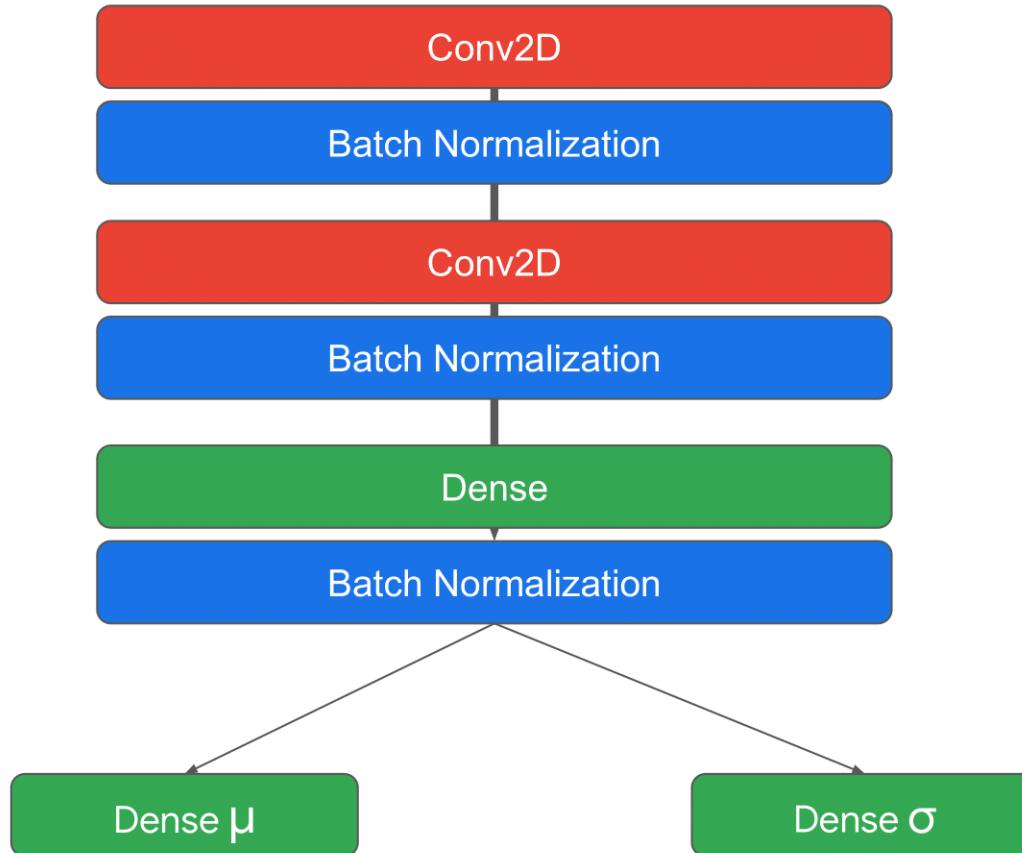
Gaussian probability density function or Normal Distribution.

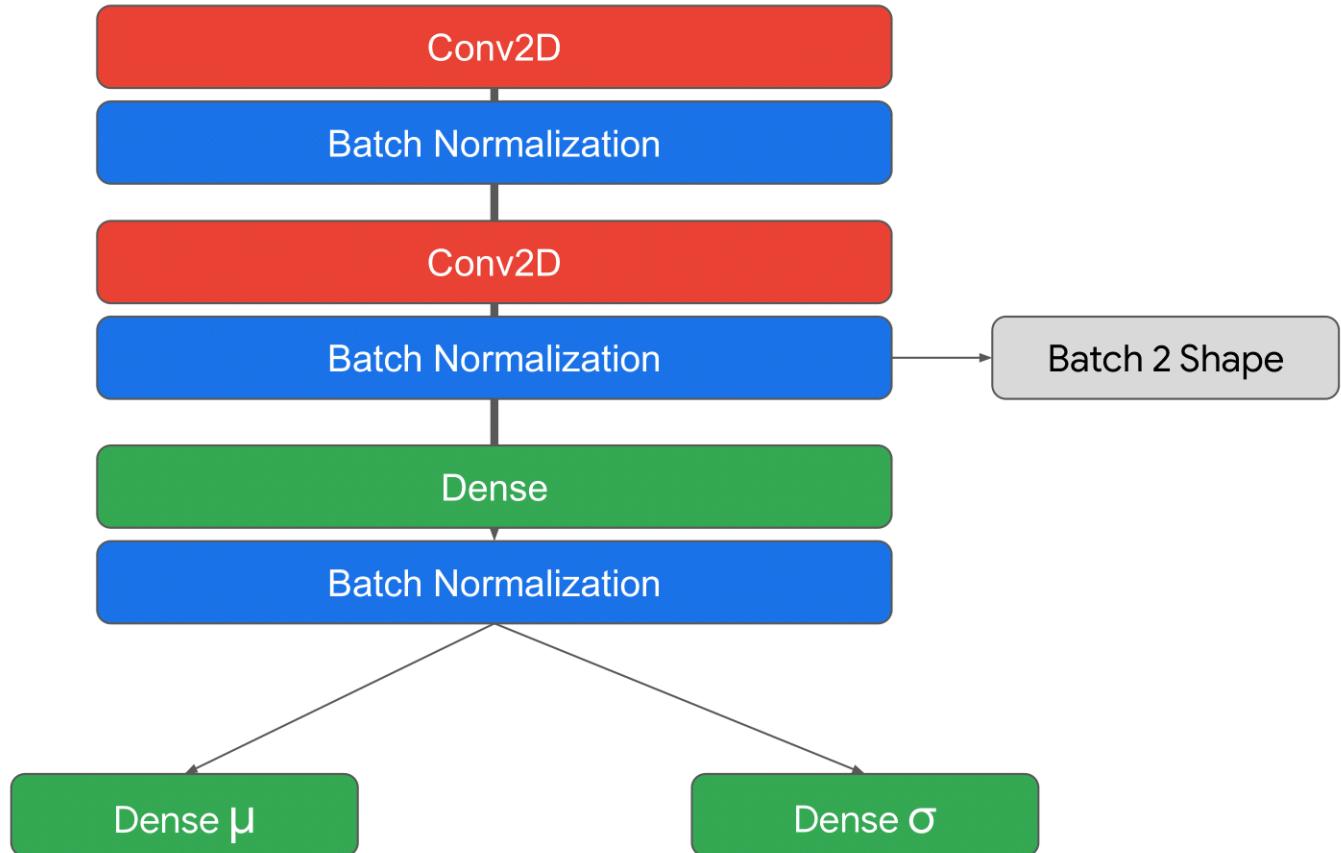
Normal Distribution is controlled by:

- μ “mean”
- σ “standard deviation”

$$N(\mu, \sigma)$$







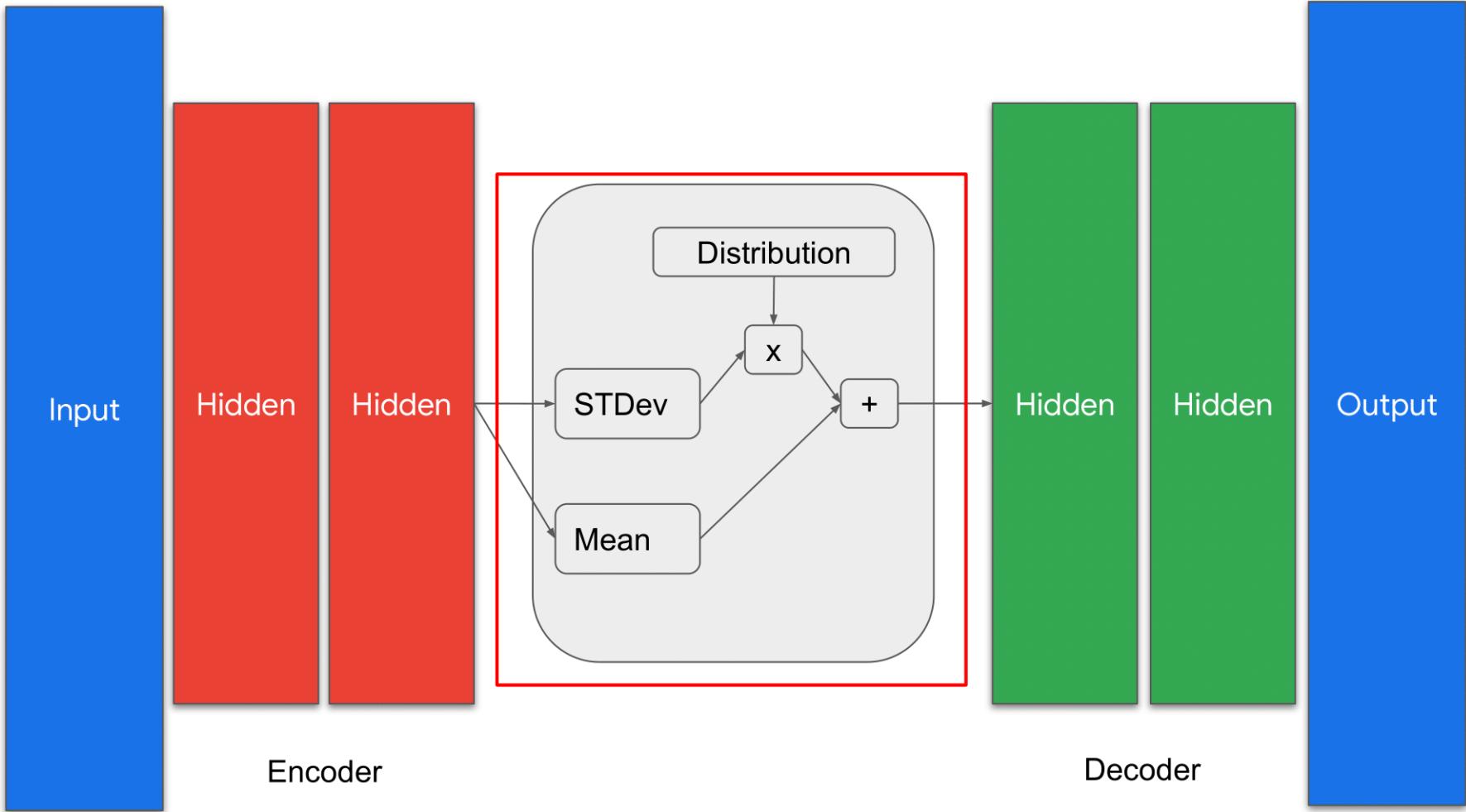
```
# This function defines the encoder's layers
def encoder_layers(inputs, latent_dim):
    x = tf.keras.layers.Conv2D(filters=32, kernel_size=3, strides=2,
                              padding="same", activation='relu',
                              name="encode_conv1")(inputs)
    x = tf.keras.layers.BatchNormalization()(x)

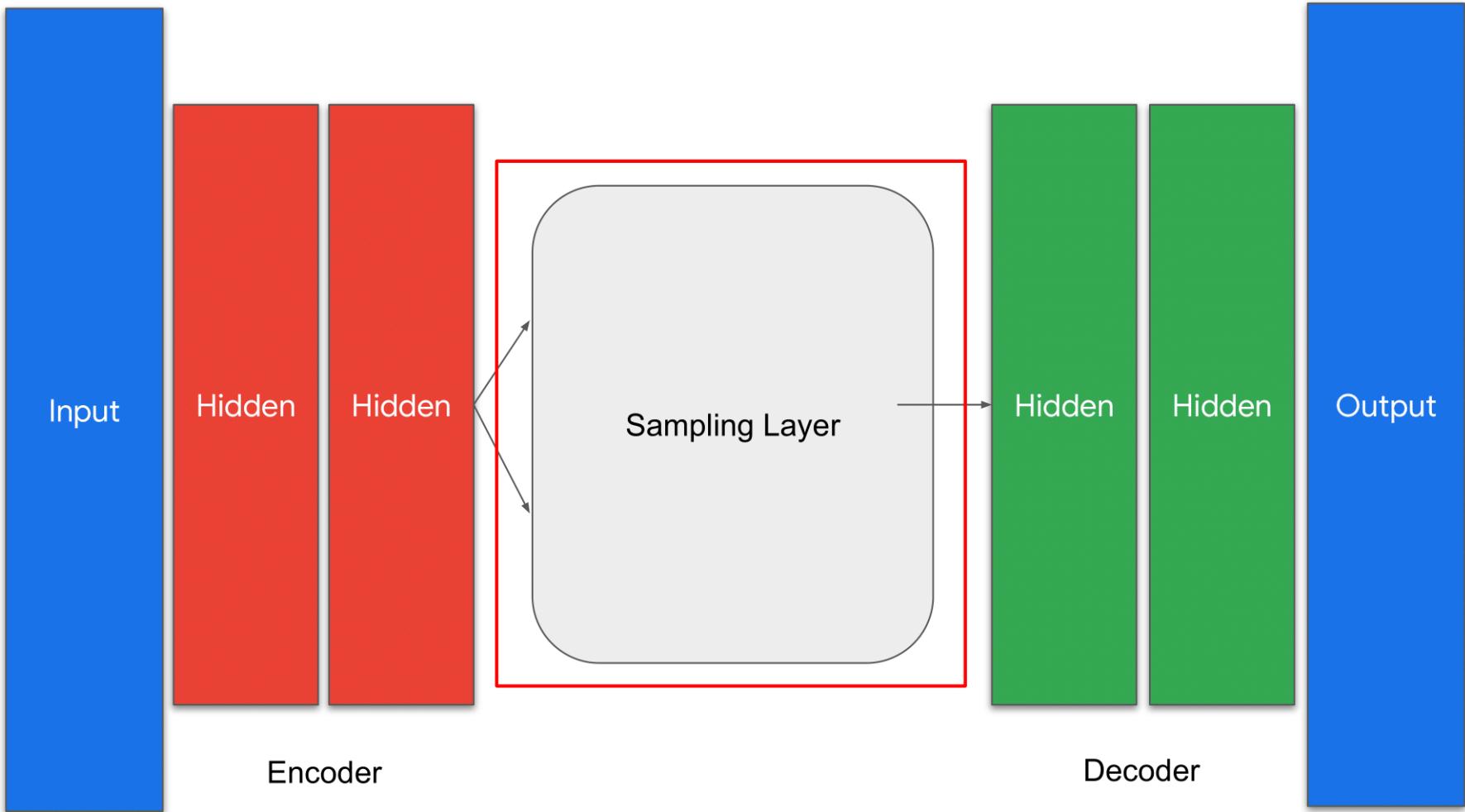
    x = tf.keras.layers.Conv2D(filters=64, kernel_size=3, strides=2,
                              padding='same', activation='relu',
                              name="encode_conv2")(x)
    batch_2 = tf.keras.layers.BatchNormalization()(x)

    x = tf.keras.layers.Flatten(name="encode_flatten")(batch_2)
    x = tf.keras.layers.Dense(20, activation='relu', name="encode_dense")(x)
    x = tf.keras.layers.BatchNormalization()(x)

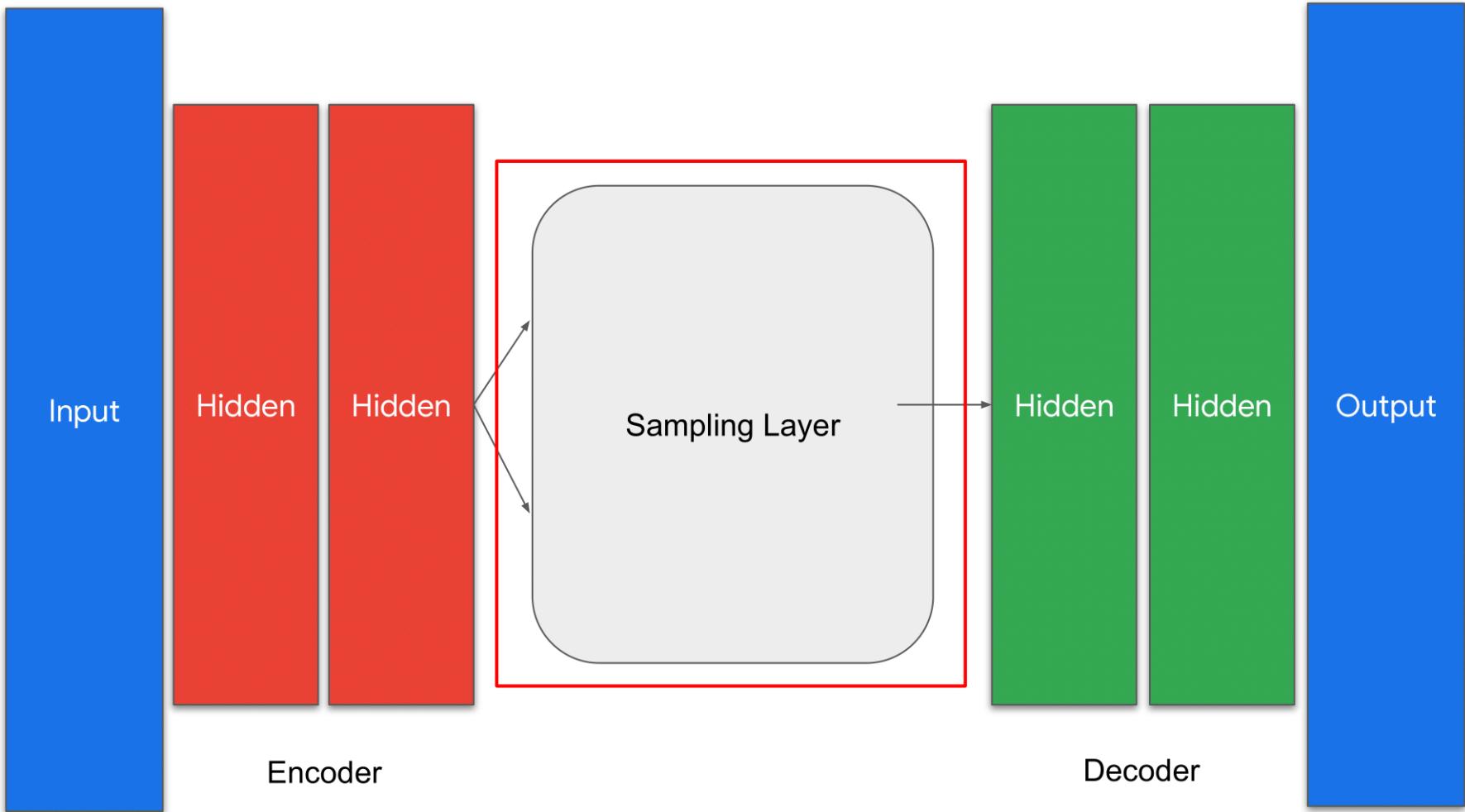
    mu = tf.keras.layers.Dense(latent_dim, name='latent_mu')(x)
    sigma = tf.keras.layers.Dense(latent_dim, name ='latent_sigma')(x)

return mu, sigma, batch_2.shape
```





```
class Sampling(tf.keras.layers.Layer):
    def call(self, inputs):
        mu, sigma = inputs
        batch = tf.shape(mu)[0]
        dim = tf.shape(mu)[1]
        epsilon = tf.keras.backend.random_normal(shape=(batch, dim))
        return mu + tf.exp(0.5 * sigma) * epsilon
```



Input

Hidden

Hidden

Sampling Layer

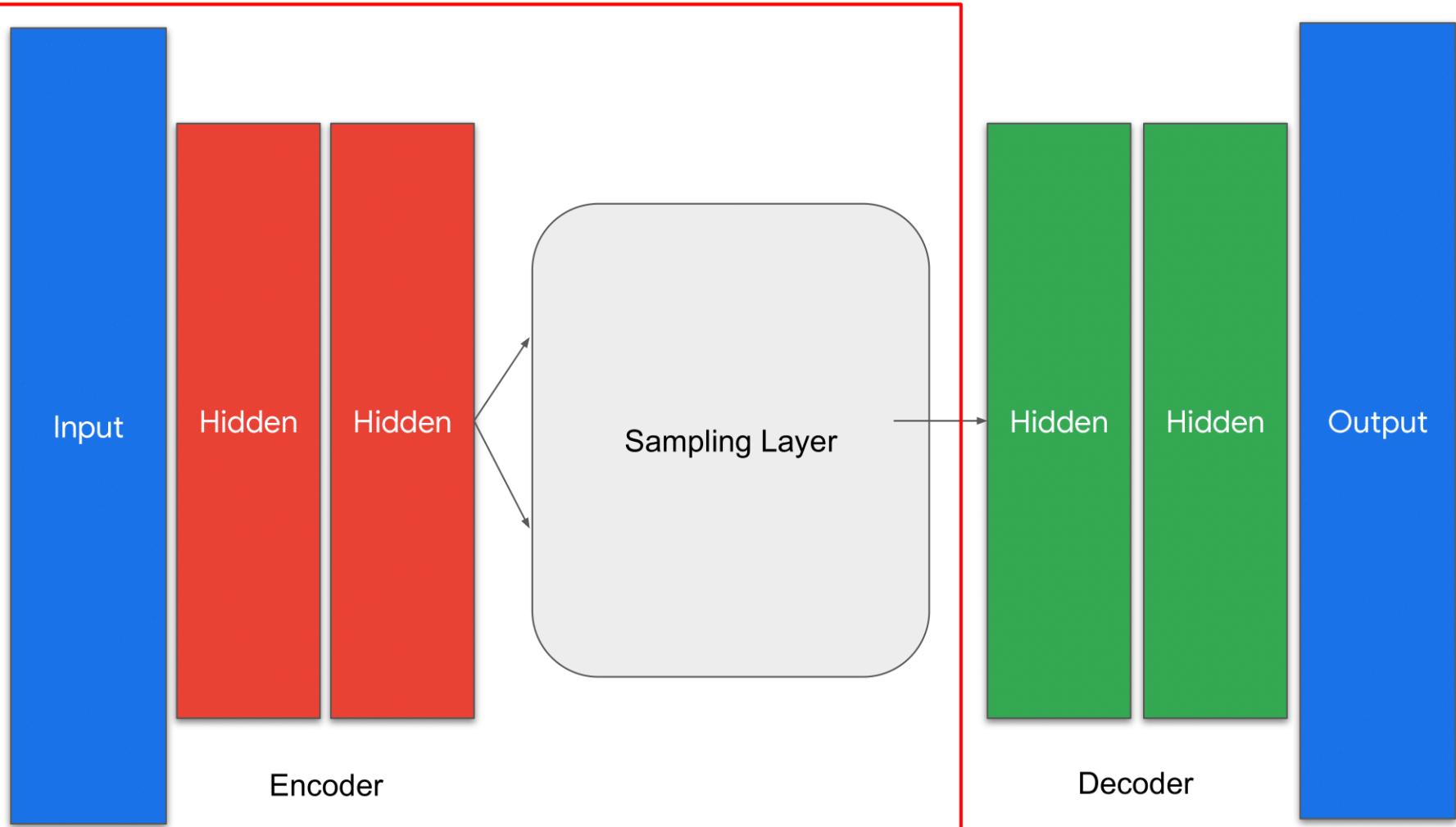
Hidden

Hidden

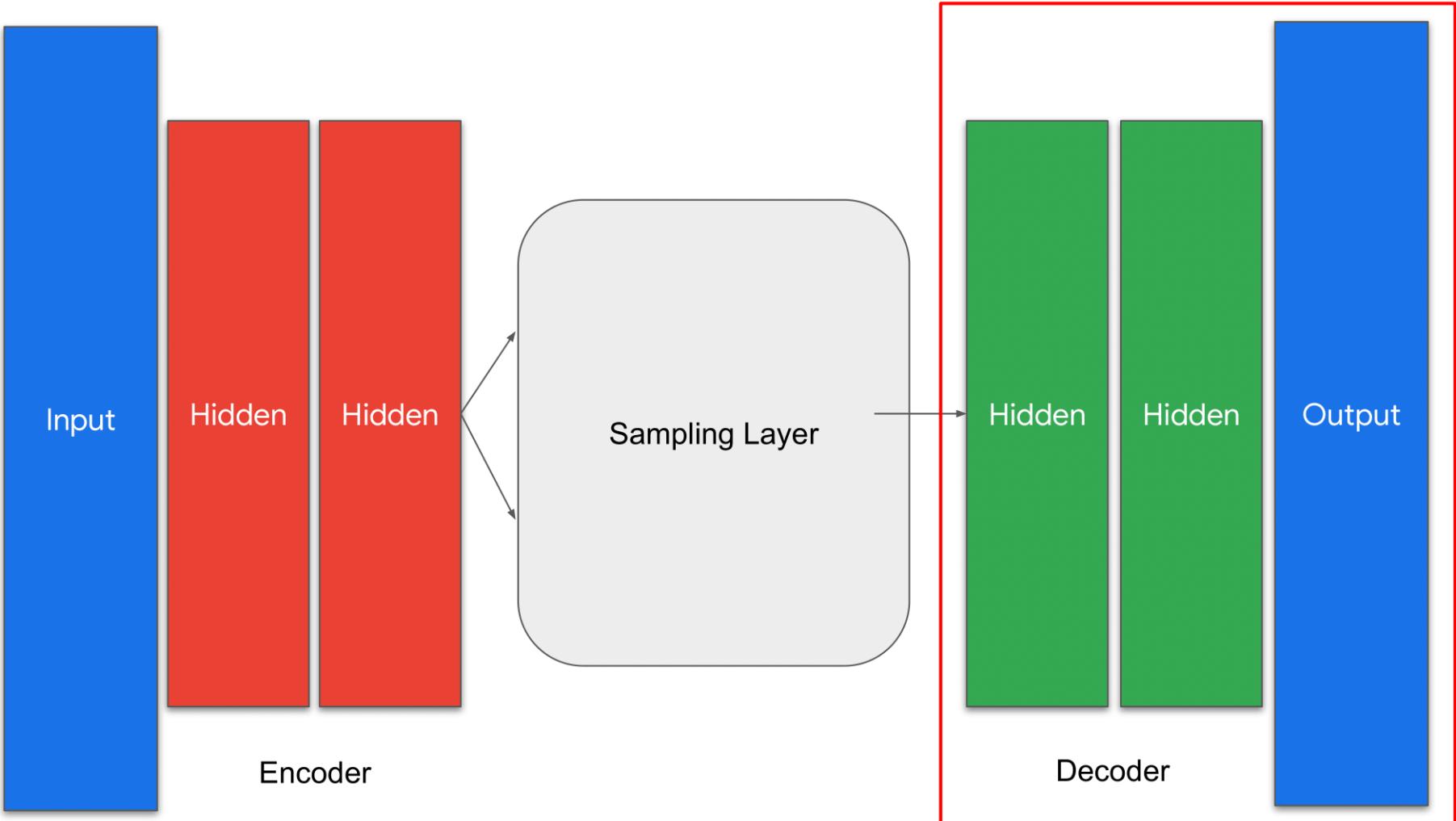
Output

Encoder

Decoder



```
def encoder_model(LATENT_DIM, input_shape):
    inputs = tf.keras.layers.Input(shape=input_shape)
    mu, sigma, conv_shape = encoder_layers(inputs, latent_dim=LATENT_DIM)
    z = Sampling()((mu, sigma))
    model = tf.keras.Model(inputs, outputs=[mu, sigma, z])
    return model, conv_shape
```



```
def decoder_layers(inputs, conv_shape):
    units = conv_shape[1] * conv_shape[2] * conv_shape[3]
    x = tf.keras.layers.Dense(units, activation = 'relu',
                             name="decode_dense1")(inputs)
    x = tf.keras.layers.BatchNormalization()(x)

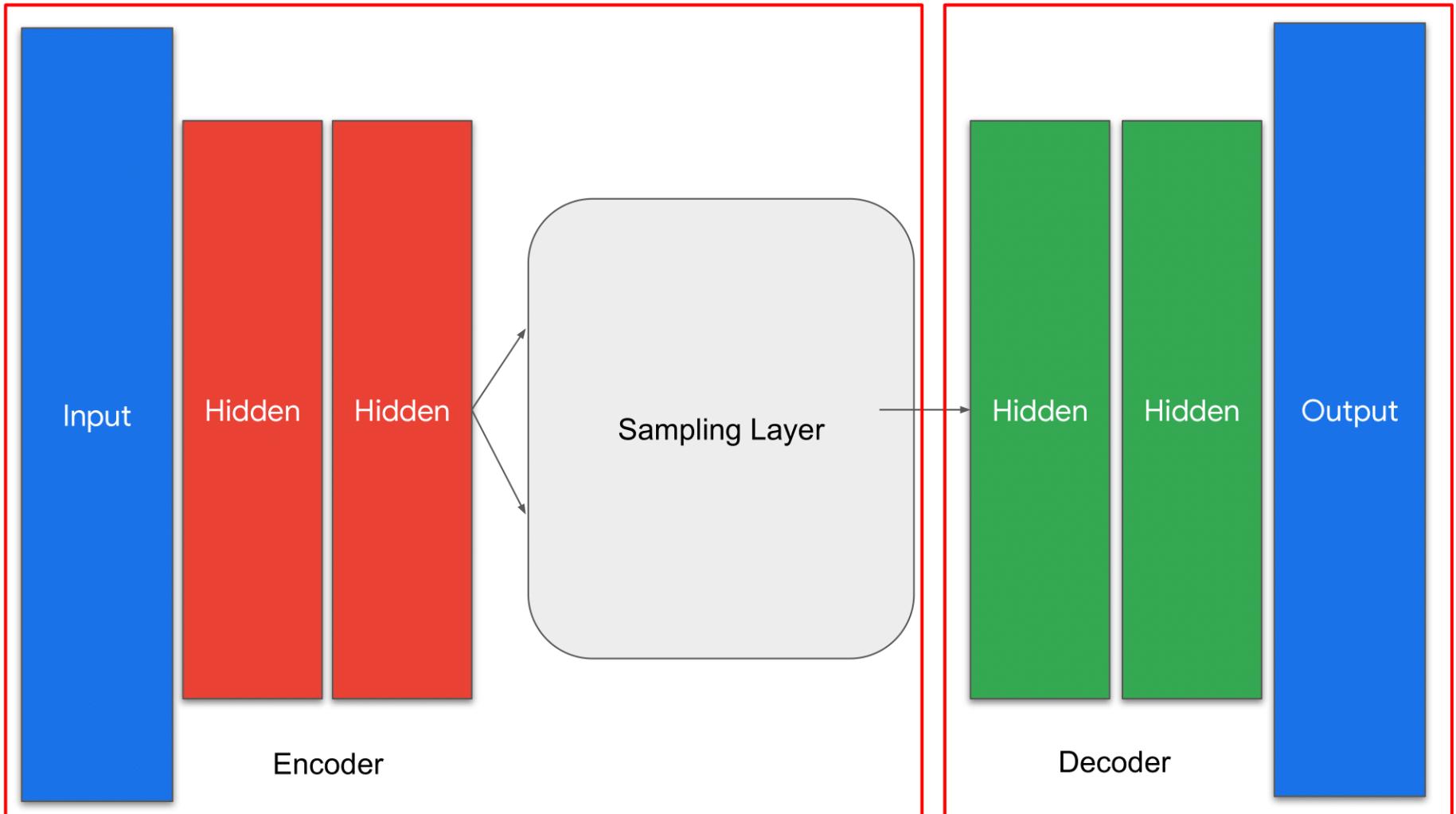
    x = tf.keras.layers.Reshape((conv_shape[1], conv_shape[2], conv_shape[3]),
                                name="decode_reshape")(x)
    x = tf.keras.layers.Conv2DTranspose(filters=64, kernel_size=3, strides=2,
                                       padding='same', activation='relu',
                                       name="decode_conv2d_2")(x)
    x = tf.keras.layers.BatchNormalization()(x)

    x = tf.keras.layers.Conv2DTranspose(filters=32, kernel_size=3, strides=2,
                                       padding='same', activation='relu',
                                       name="decode_conv2d3")(x)
    x = tf.keras.layers.BatchNormalization()(x)

    x = tf.keras.layers.Conv2DTranspose(filters=1, kernel_size=3, strides=1, padding='same',
                                       activation='sigmoid', name="decode_final")(x)

return x
```

```
def decoder_model(latent_dim, conv_shape):
    inputs = tf.keras.layers.Input(shape=(latent_dim,))
    outputs = decoder_layers(inputs, conv_shape)
    model = tf.keras.Model(inputs, outputs)
    return model
```



```
# Define a kl reconstruction loss function

def kl_reconstruction_loss(inputs, outputs, mu, sigma):
    kl_loss = 1 + sigma - tf.square(mu) - tf.math.exp(sigma)
    return tf.reduce_mean(kl_loss) * -0.5
```

```
def vae_model(encoder, decoder, input_shape):
    inputs = tf.keras.layers.Input(shape=input_shape)
    mu = encoder(inputs)[0]
    sigma = encoder(inputs)[1]
    z = encoder(inputs)[2]
    reconstructed = decoder(z)
    model = tf.keras.Model(inputs=inputs, outputs=reconstructed)
    loss = kl_reconstruction_loss(inputs, z, mu, sigma)
    model.add_loss(loss)
    return model
```

```
for epoch in range(epochs):
    for step, x_batch_train in enumerate(train_dataset):
        with tf.GradientTape() as tape:
            reconstructed = vae(x_batch_train)
            flattened_inputs = tf.reshape(x_batch_train, shape=[-1])
            flattened_outputs = tf.reshape(reconstructed, shape=[-1])
            loss = bce_loss(flattened_inputs, flattened_outputs) * 784
            loss += sum(vae.losses) # Add KLD regularization loss

        grads = tape.gradient(loss, vae.trainable_weights)
        optimizer.apply_gradients(zip(grads, vae.trainable_weights))
```

epoch: 99, step: 400



Variational Autoencoders

by Paul Hand
Northeastern University

Outline:

Generative Models and Autoencoders

Variational Lower Bound (VLB)

Optimizing VLB + Variational Autoencoders

Resource: Kingma and Welling 2019, Chapter 2,
"Introduction to Variational Autoencoders"

Generative Models

A model that can sample from a learned distribution

8 6 / 7 8 1 4 8 2 8
9 6 8 3 9 6 8 3 1 9
3 3 9 1 3 6 9 1 7 9
8 9 0 8 6 9 1 9 6 3
8 2 3 3 3 3 1 3 8 6
6 4 4 8 6 1 6 6 6 6
9 5 2 6 6 5 1 8 9 9
9 9 8 9 3 1 2 8 2 3
0 4 6 1 2 3 2 0 8 8
9 7 5 4 9 3 4 8 5 1

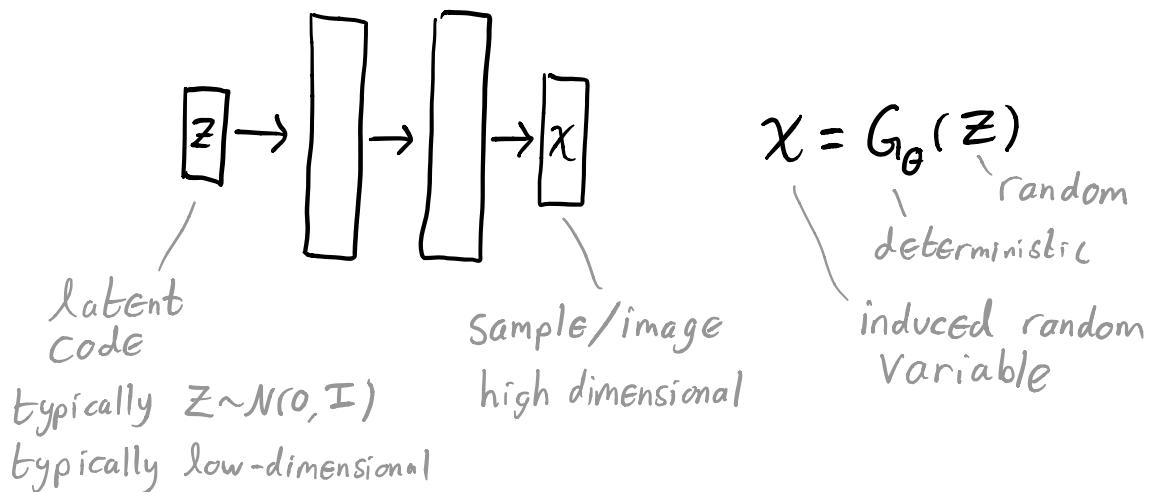
(a) 2-D latent space

(Kingma and Welling, 2014)



(Razavi et al. 2019)

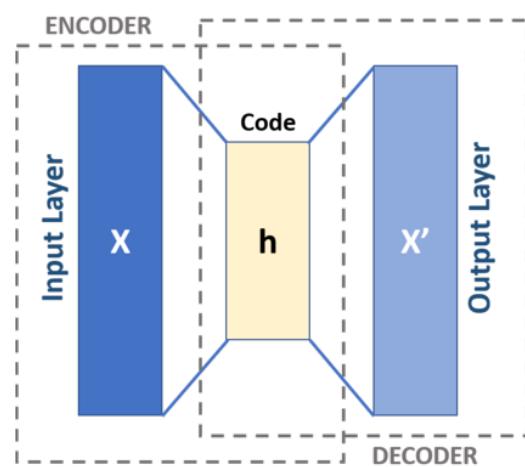
In some generative models, samples are generated by a net applied to a random latent code



Autoencoders

Autoencoders attempt to reconstruct input signals/images by learning mappings to and from a code

$$\text{Want}^{\circ}: X' = D_\theta(E_\phi(x)) \approx x$$



$$\min_{\theta, \phi} \sum_{i=1}^n \|D_\theta(E_\phi(x_i)) - x_i\|^2 \quad \text{w/ } \{x_i\}_{i=1...n} \text{ is dataset}$$

A (plain) autoencoder is not a generative model as it does not define a distribution

Training a low latent-dimensional generative model by likelihood

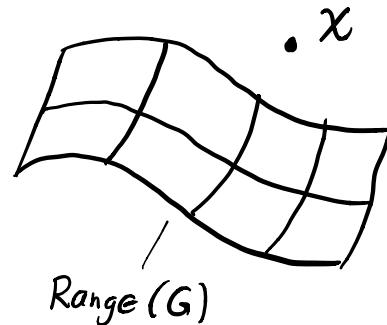
Given data $\{x_i\}_{i=1...n}$, train a gen. model to maximize the likelihood of the observed data

If gen. model

$$G_\theta : \mathbb{R}^k \rightarrow \mathbb{R}^d \quad \text{w/ } k < d, \\ z \mapsto x$$

then $p(x)=0$ almost everywhere

So, can't directly optimize likelihood



To have nonzero likelihood everywhere, defining noisy observation model

$$P_\theta(x|z) = N(x; G_\theta(z), \gamma I)$$

Under a simple prior $p(z)$, this induces

a joint distribution $P_\theta(x, z)$

Now $p(x) = \int p(z) p(x|z) dz$

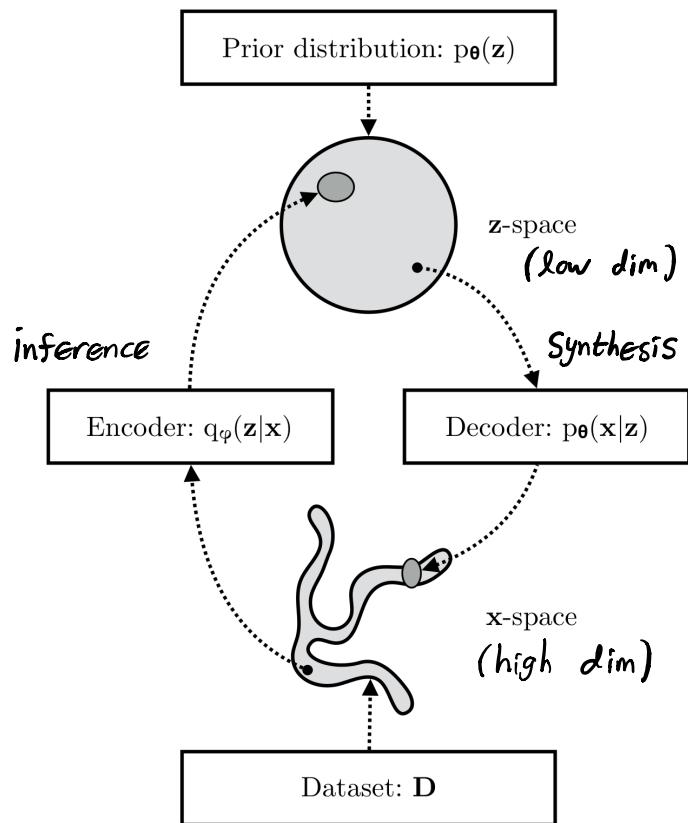
/
intractable to evaluate at each iteration
optimize a lower bound instead

Variational Lower Bound

Setup:

Data generated by
 $z \sim p(z)$ prior
 $x \sim P_\theta(x|z)$

Use $q_\varphi(z|x)$ as
proxy for $P_\theta(z|x)$



(Kingma and Welling 2019)

Find \circlearrowleft lower bound to $P_\theta(x)$

$$\begin{aligned}
 \log P_\theta(x) &= \mathbb{E}_{z \sim q_\phi(z|x)} \log P_\theta(x) = \mathbb{E}_{z \sim q_\phi(z|x)} \log \frac{P_\theta(x, z)}{P_\theta(z|x)} \\
 &= \mathbb{E}_{z \sim q_\phi(z|x)} \log \left(\frac{P_\theta(x, z)}{q_\phi(z|x)} \cdot \frac{q_\phi(z|x)}{P_\theta(z|x)} \right) \\
 &= \underbrace{\mathbb{E}_{z \sim q_\phi(z|x)} \log \frac{P_\theta(x, z)}{q_\phi(z|x)}}_{\mathcal{L}_{\theta, \phi}(x)} + \underbrace{\mathbb{E}_{z \sim q_\phi(z|x)} \log \frac{q_\phi(z|x)}{P_\theta(z|x)}}_{D_{KL}(q_\phi(z|x) \parallel P_\theta(z|x))}
 \end{aligned}$$

Variational Lower Bound (VLB)

Evidence lower bound (ELBO)

Define \circlearrowleft $D_{KL}(q||P) = \mathbb{E}_{z \sim q} \log \frac{q(z)}{P(z)}$

- Note \circlearrowleft
- $D_{KL}(q||P) \neq D_{KL}(P||q)$
 - Measure of how far P is from q
 - $D_{KL}(q||P) \geq 0$ (and is 0 iff $P=q$)

So, $\log P_\theta(x) \geq \mathbb{E}_{z \sim q_\phi(z|x)} \log \frac{P_\theta(x, z)}{q_\phi(z|x)} = \mathcal{L}_{\theta, \phi}(x)$

Interpretation of VLB

$$\begin{aligned}
 L_{\theta, \psi}(x) &= \mathbb{E}_{z \sim q_\psi(z|x)} \log \frac{P_\theta(x|z)}{q_\psi(z|x)} = \mathbb{E}_{z \sim q_\psi(z|x)} \log P_\theta(x|z) \frac{P(z)}{q_\psi(z|x)} \\
 &= \underbrace{\mathbb{E}_{z \sim q_\psi(z|x)} \log P_\theta(x|z)}_{\text{reconstruction error}} + \underbrace{\mathbb{E}_{z \sim q_\psi(z|x)} \log \frac{P(z)}{q_\psi(z|x)}}_{\text{regularization}}
 \end{aligned}$$

First term: $P_\theta(x|z) = N(x; G_\theta(z), \gamma I)$
 $\Rightarrow \log P_\theta(x|z) = -\frac{1}{2\gamma} \|x - G_\theta(z)\|^2 + \text{constant}$

So $\mathbb{E}_{z \sim q_\psi} \log P_\theta(x|z)$ is expected ℓ_2 reconstruction error under the encoder model

Maximizing VLB encourages q_ψ to be point mass

Second term: $\mathbb{E}_{z \sim q_\psi(z|x)} \log \frac{P(z)}{q_\psi(z|x)} = -D_{KL}(q_\psi(z|x) || P(z))$

So maximizing VLB $L_{\theta, \psi}$ pushes $q_\psi(z|x)$ toward $P(z)$. Prevents q_ψ from being a point mass.

Makes $q_\psi(z|x)$ more like standard normal

Maximizing VLB $\mathcal{L}_{\theta, \varphi}$:

- Roughly maximizes $P(x)$
- Minimizes KL divergence of $q_{\varphi}(z|x)$ and $P_{\theta}(z|x)$, making q_{φ} better

Main

Idea: Instead of optimizing $\sum_{i=1}^n \log P_{\theta}(x_i)$,
optimize $\sum_{i=1}^n \mathcal{L}_{\theta, \varphi}(x_i)$
w/ $\mathcal{L}_{\theta, \varphi}(x) = \mathbb{E}_{z \sim q_{\varphi}(z|x)} \log \frac{P_{\theta}(x, z)}{q_{\varphi}(z|x)}$

Optimizing Variational Lower Bound

$$\max_{\theta, \varphi} \sum_{i=1}^n \mathcal{L}_{\theta, \varphi}(x_i)$$

One possibility:

for each x_i , find best
 $q_{\varphi}(z|x_i)$ by multiple gradient
steps in φ . Then gradient ascend
in θ .

Expensive inference updates

Instead:

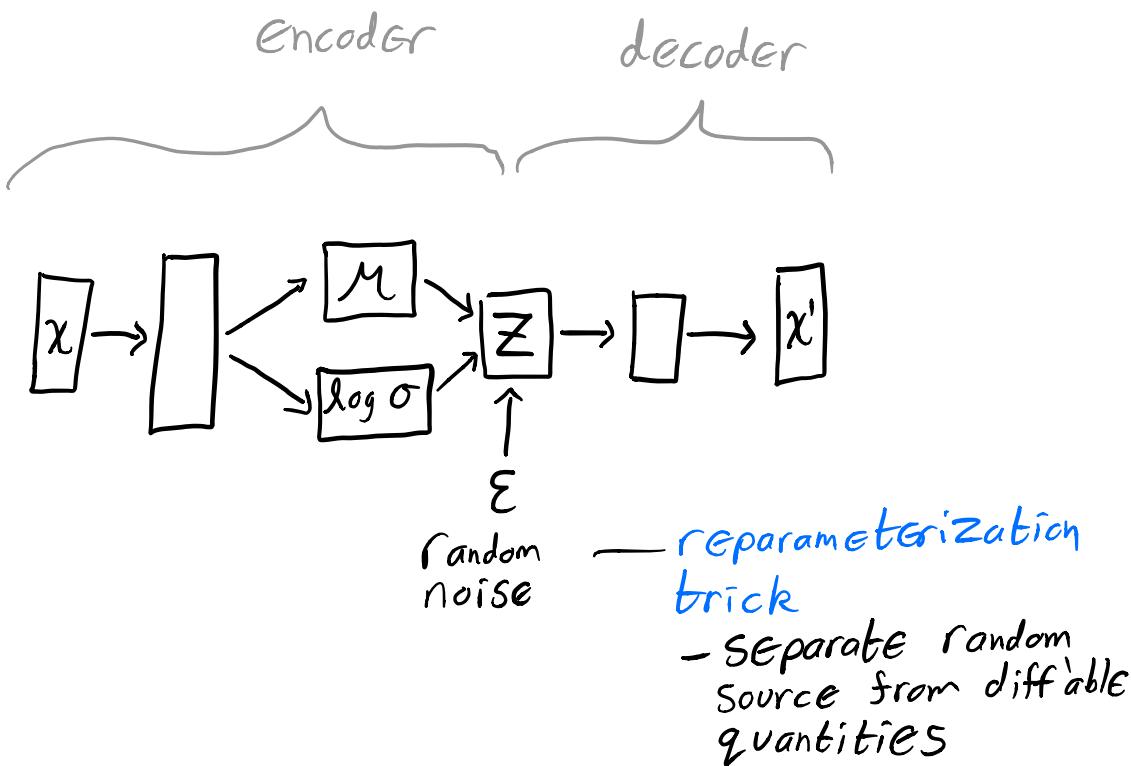
amortize the inference costs
by learning an inference net

$$x \mapsto (\mu, \Sigma) \text{ w/ } q_{\phi}(z|x) = N(z; \mu(x), \Sigma(x))$$

or $\sigma(x)I$

Parameters of inference model
are shared between data points

Variational Autoencoder architecture



Stochastic Gradient Optimization of VLB

Dataset $D = \{x_i\}_{i=1 \dots n}$

Solve $\max_{\theta, \psi} \sum_{x_i \in D} \mathcal{L}_{\theta, \psi}(x_i)$

$$\text{w/ } \mathcal{L}_{\theta, \psi}(x) = \mathbb{E}_{z \sim q_\psi(z|x)} \log \frac{P_\theta(x, z)}{q_\psi(z|x)}$$

Computing $\nabla_{\theta, \psi} \mathcal{L}_{\theta, \psi}(x_i)$ is intractable,
but there are unbiased estimators

Easy to get unbiased $\nabla_\theta \mathcal{L}_{\theta, \psi} \circledast$:

$$\begin{aligned} \nabla_\theta \mathcal{L}_{\theta, \psi}(x) &= \nabla_\theta \mathbb{E}_{z \sim q_\psi(z|x)} [\log P_\theta(x, z) - \log q_\psi(z|x)] \\ &= \mathbb{E}_{z \sim q_\psi} \nabla_\theta \log P_\theta(x, z) \\ &\stackrel{\sim}{=} \nabla_\theta \log P_\theta(x, z) \text{ w/ } z \sim q_\psi(z|x) \\ &\quad \text{unbiased estimate} \end{aligned}$$

Not as easy to get unbiased $\nabla_{\varphi} \mathcal{L}_{\theta, \varphi}$:

$$\begin{aligned}\nabla_{\varphi} \mathcal{L}_{\theta, \varphi}(x) &= \nabla_{\varphi} \mathbb{E}_{z \sim q_{\varphi}(z|x)} [\log P_{\theta}(x, z) - \log q_{\varphi}(z|x)] \\ &\neq \mathbb{E}_{z \sim q_{\varphi}} [\nabla_{\varphi} (\log P_{\theta}(x, z) - \log q_{\varphi}(z|x))]\end{aligned}$$

Recall, $q_{\varphi}(z|x) = N(z; \mu(x), \sigma(x)\mathbf{I})$
 $= \mu(x) + \sigma(x) \cdot \varepsilon, \text{ w/ } \varepsilon \sim N(0, \mathbf{I})$

$$\begin{aligned}\mathcal{L}_{\theta, \varphi}(x) &= \mathbb{E}_{z \sim q_{\varphi}(z|x)} (\log P_{\theta}(x, z) - \log q_{\varphi}(z|x)) \\ &= \mathbb{E}_{\varepsilon \sim P(\varepsilon)} (\log P_{\theta}(x, z) - \log q_{\varphi}(z|x))\end{aligned}$$

Form estimator of $\mathcal{L}_{\theta, \varphi}(x)$ as $\hat{\mathcal{L}}_{\theta, \varphi}(x)$ by:

$$\varepsilon \sim P(\varepsilon)$$

$$z = \mu_{\varphi}(x) + \sigma_{\varphi}(x) \varepsilon = g(\varphi, x, \varepsilon)$$

$$\hat{\mathcal{L}}_{\theta, \varphi}(x) = \log P_{\theta}(x, z) - \log q_{\varphi}(z|x)$$

Unbiased estimate of $\nabla_{\varphi} \mathcal{L}_{\theta, \varphi}(x)$

$$\nabla_{\varphi} \hat{\mathcal{L}}_{\theta, \varphi}(x)$$

Note: $\mathbb{E}_{\varepsilon \sim p(\varepsilon)} \hat{\mathcal{L}}_{\theta, \varphi}(x) = \mathcal{L}_{\theta, \varphi}(x)$

So $\mathbb{E}_{\varepsilon \sim p(\varepsilon)} \nabla_{\varphi} \hat{\mathcal{L}}_{\theta, \varphi}(x) = \nabla_{\varphi} \mathbb{E}_{\varepsilon \sim p(\varepsilon)} \hat{\mathcal{L}}_{\theta, \varphi} = \nabla_{\varphi} \mathcal{L}_{\theta, \varphi}$

Optimize VAE parameters w/ stochastic gradients.

Can extend models to be more sophisticated,
e.g. $\Sigma(x)$ vs $\sigma I(x)$ as inference model.

Key points:

- optimize a lower bound to likelihood
- lower bound has terms for reconstruction and for regularization
- maintain an inference model for $Z|x$ in place of intractable true distribution

- reparameterization trick allows backpropagating on mean and variance of inference model
- VAEs have been trained with photorealistic outputs