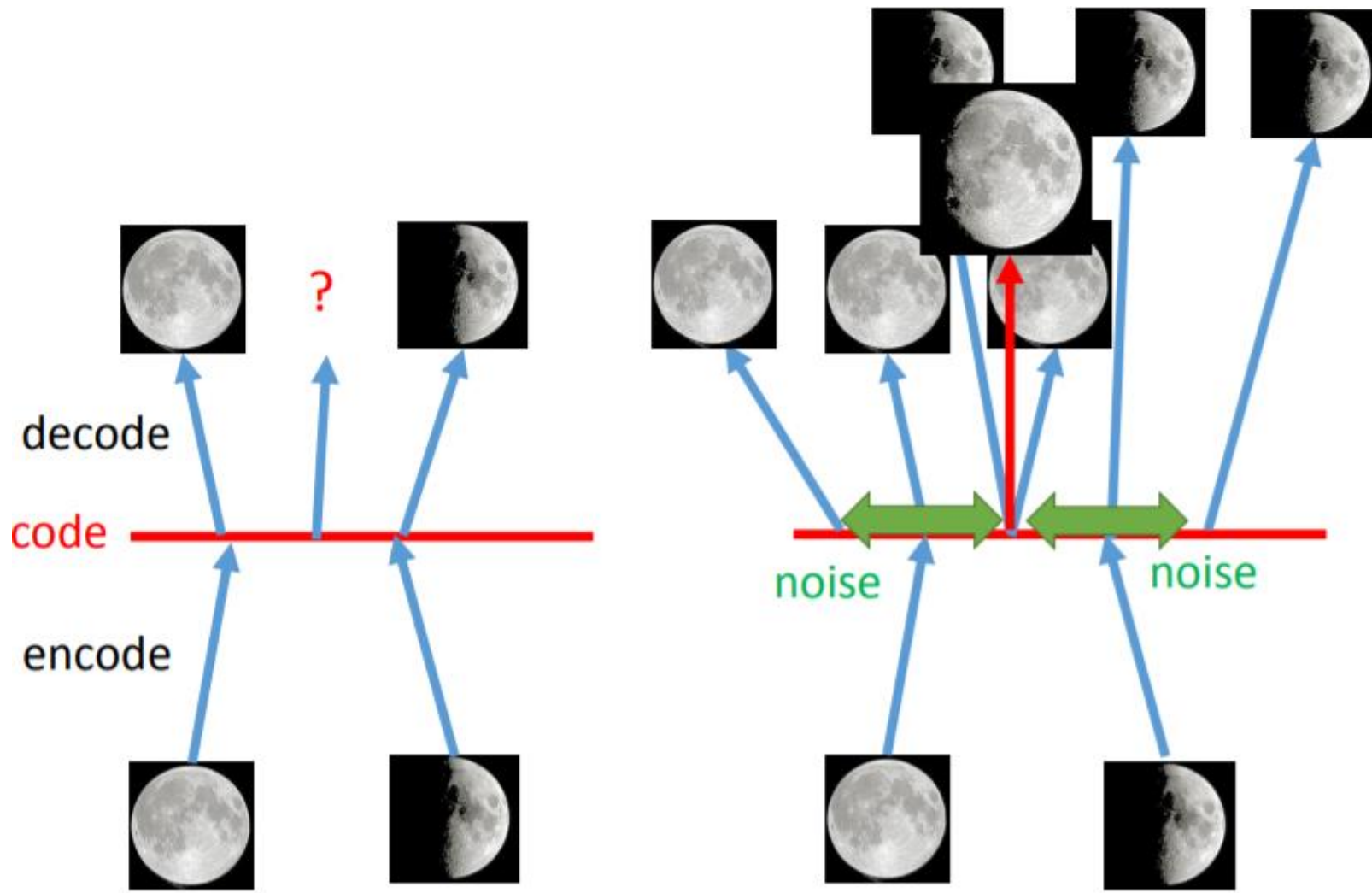# Vibrational Auto-Encoder (VAE)

# Why VAE?

Assume 1-d code
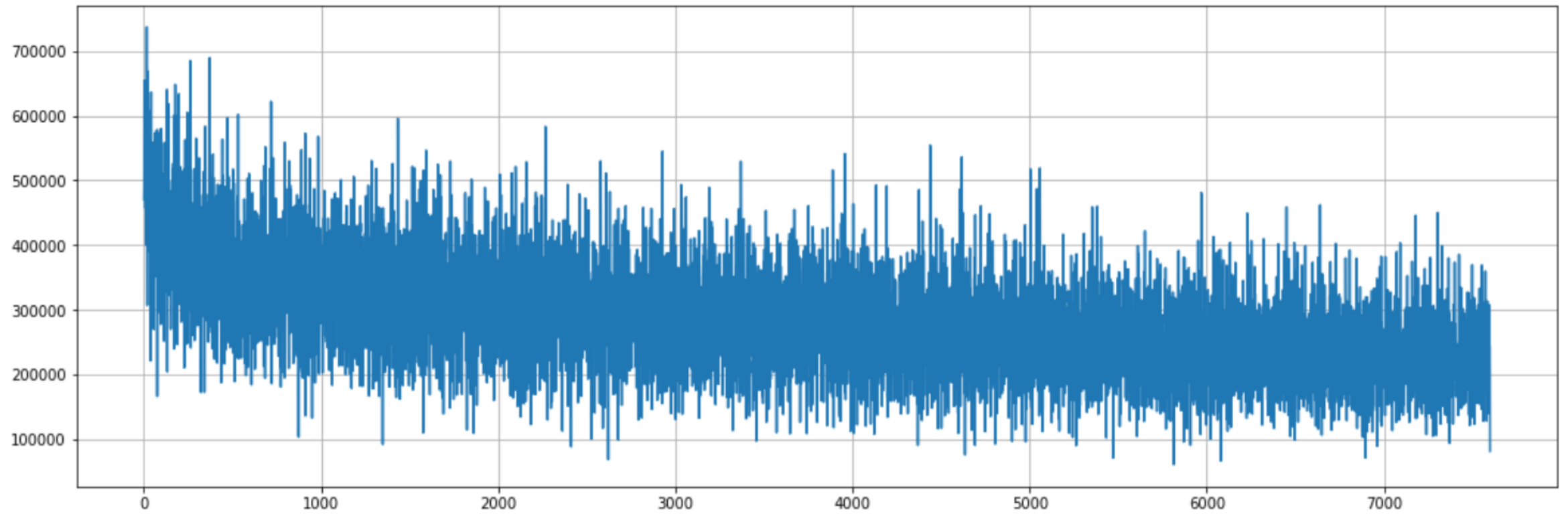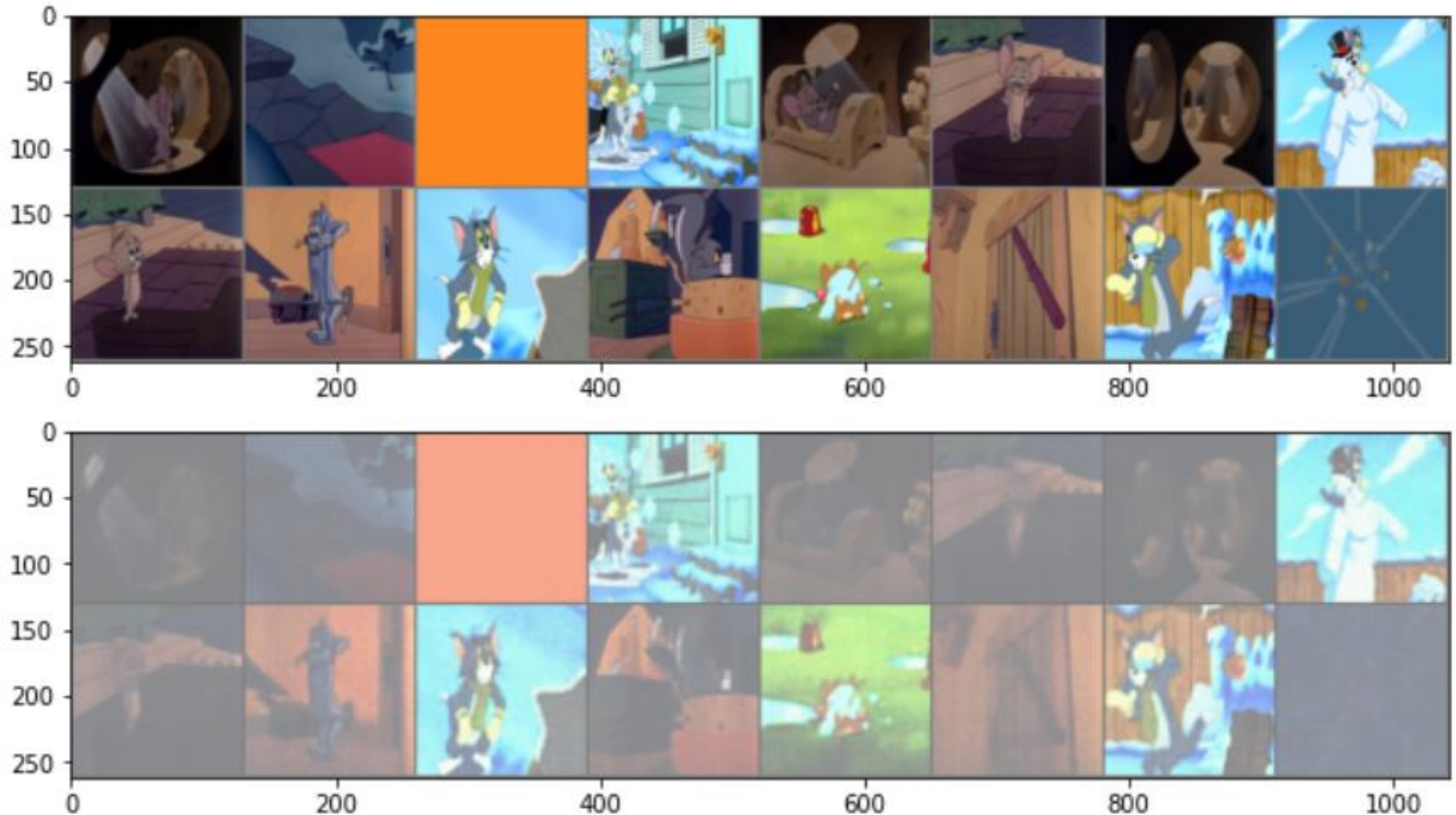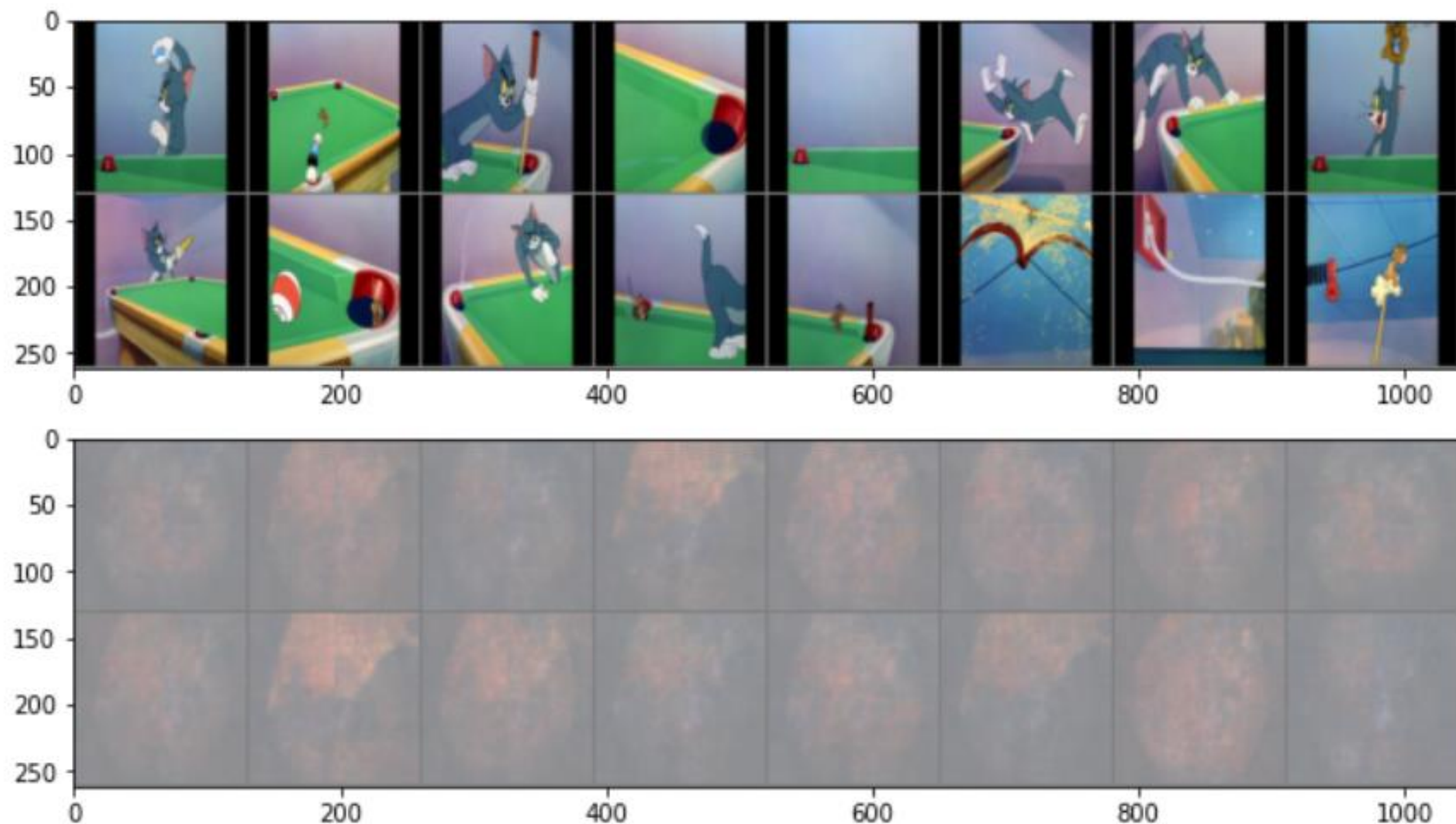
# Practice

- Run "7.2.Conv_VAE.ipynb"

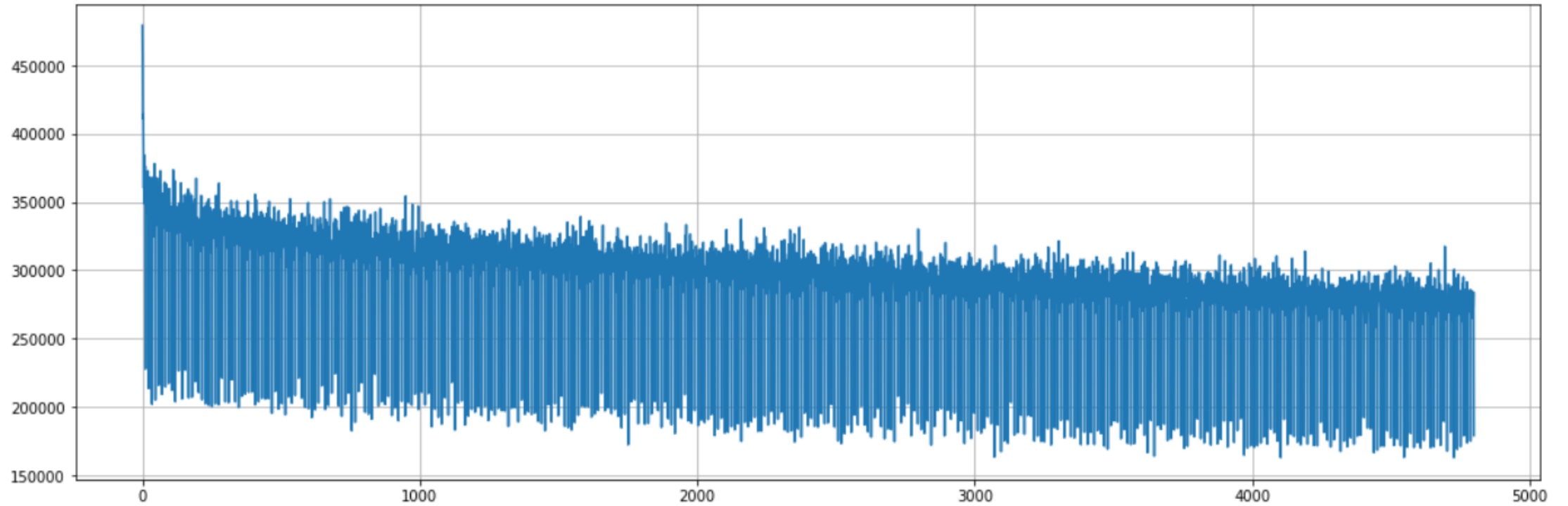# Loss plot of epoch 0-400

# Training images recovered after training for 400 epochs

# NN fails to recover un-seen test images if trained for 400 epochs

# Loss plot of epoch 400-800

# Training images recovered after training for 800 epochs

# After training for 800 epochs, the NN can recover un-seen test images
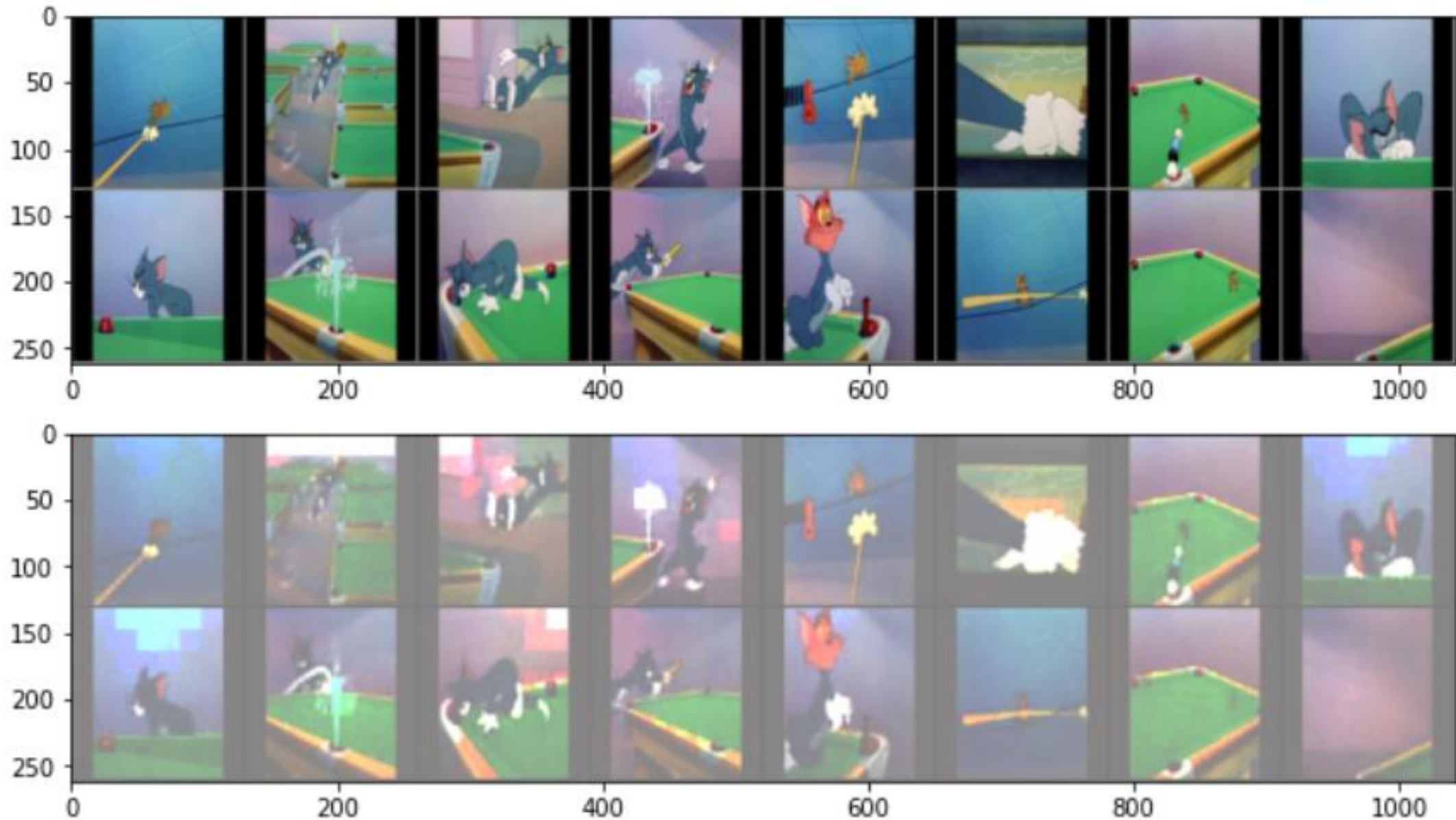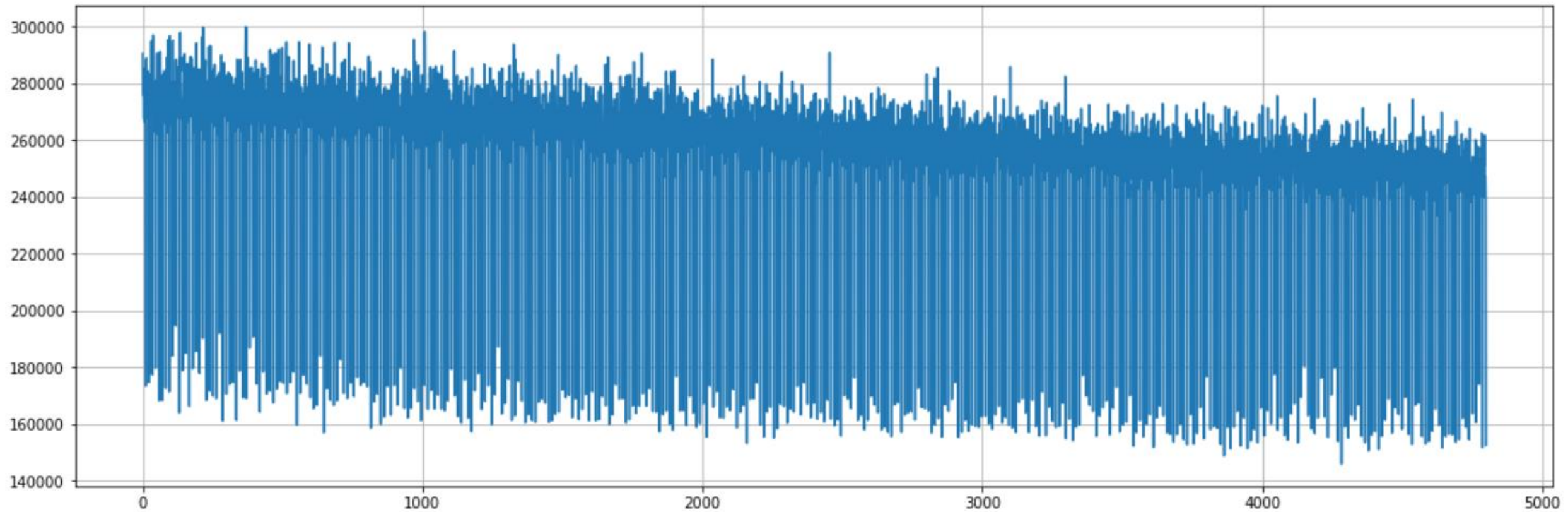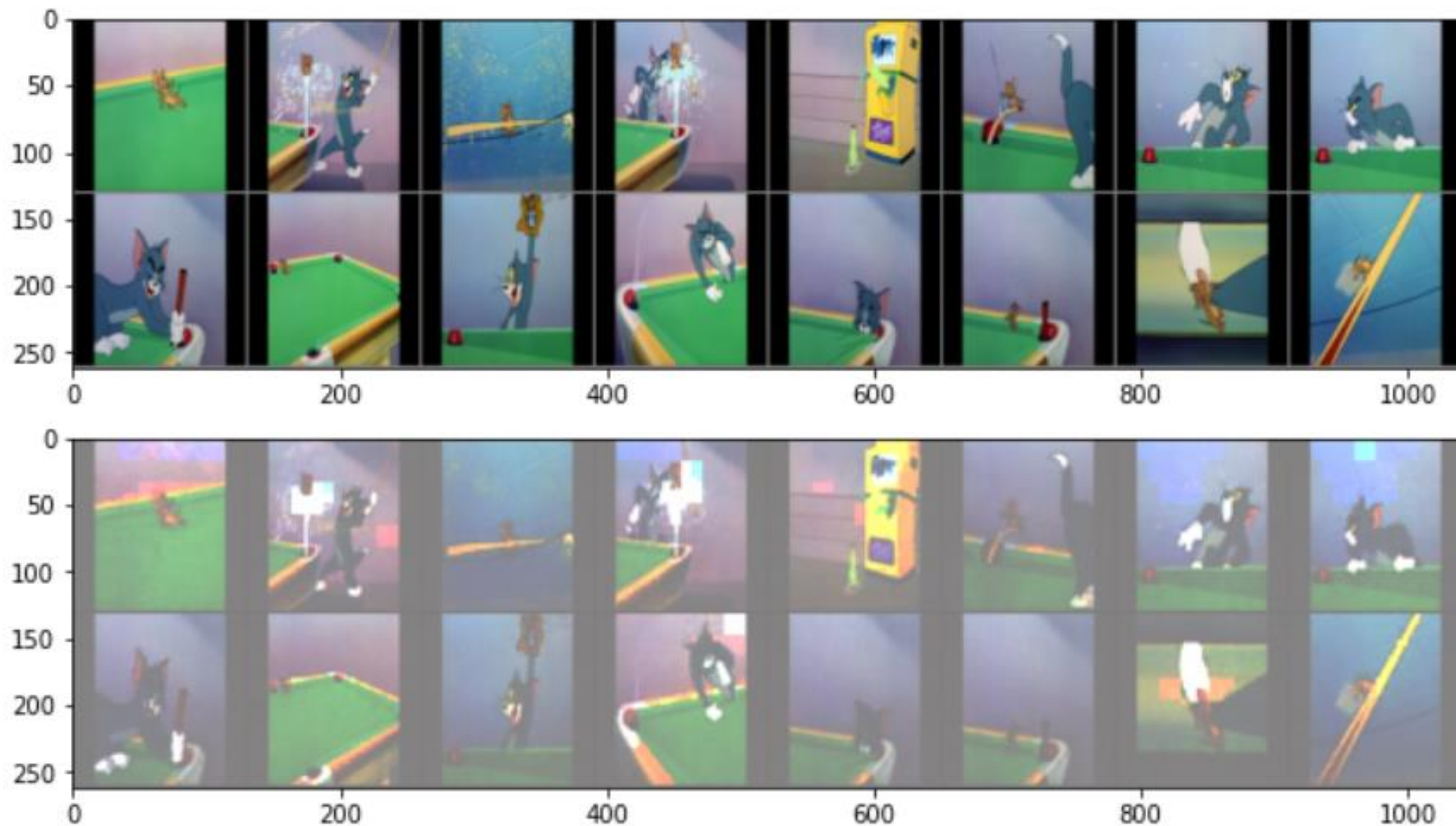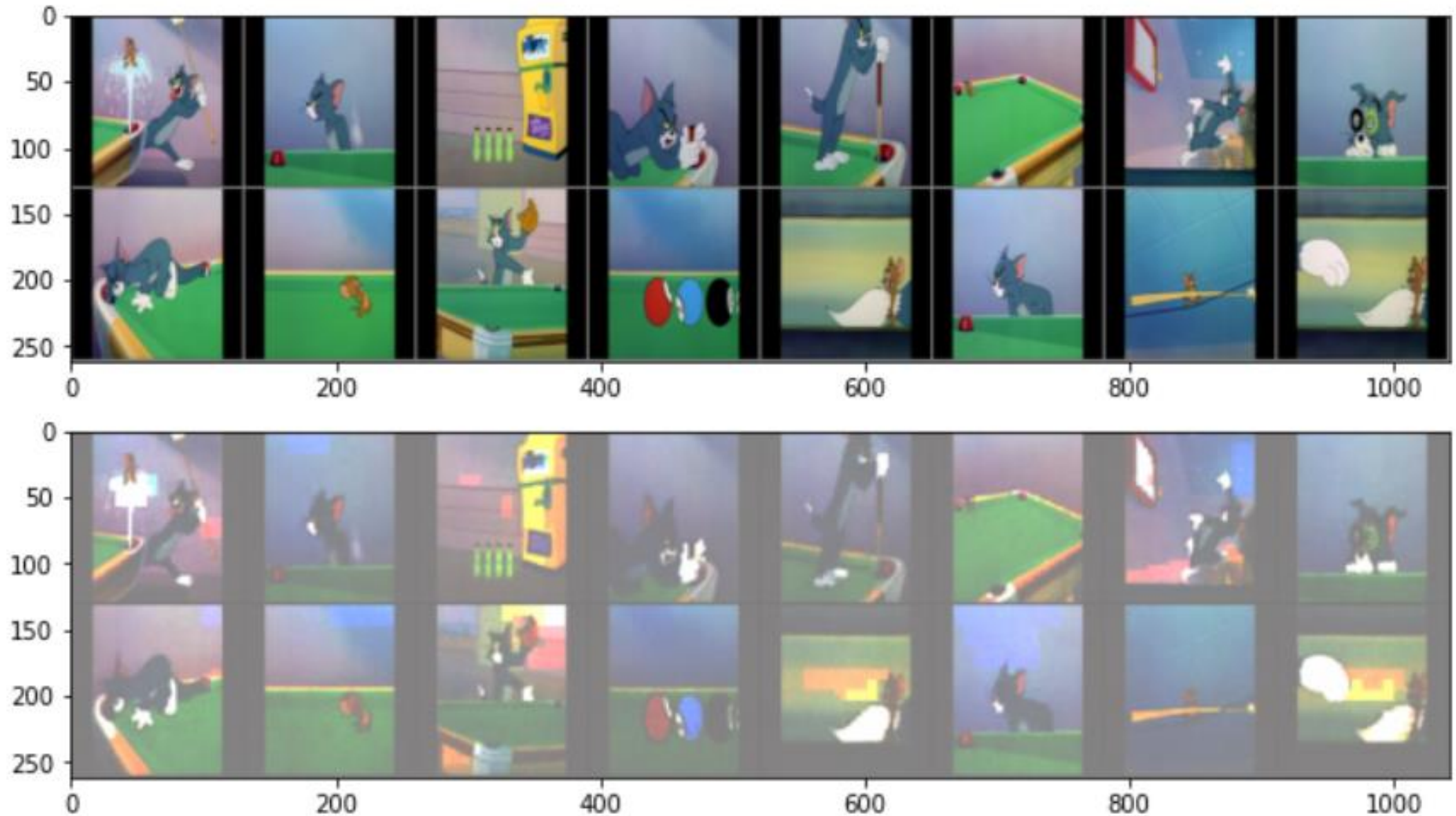
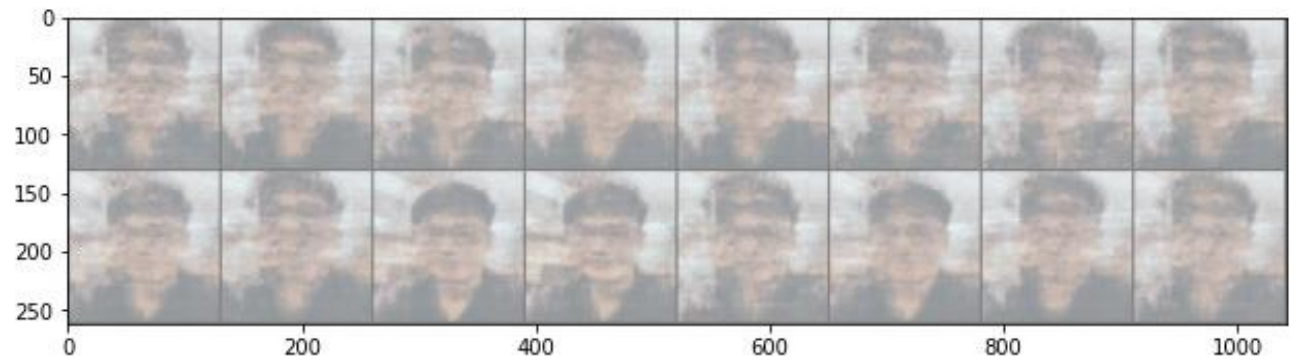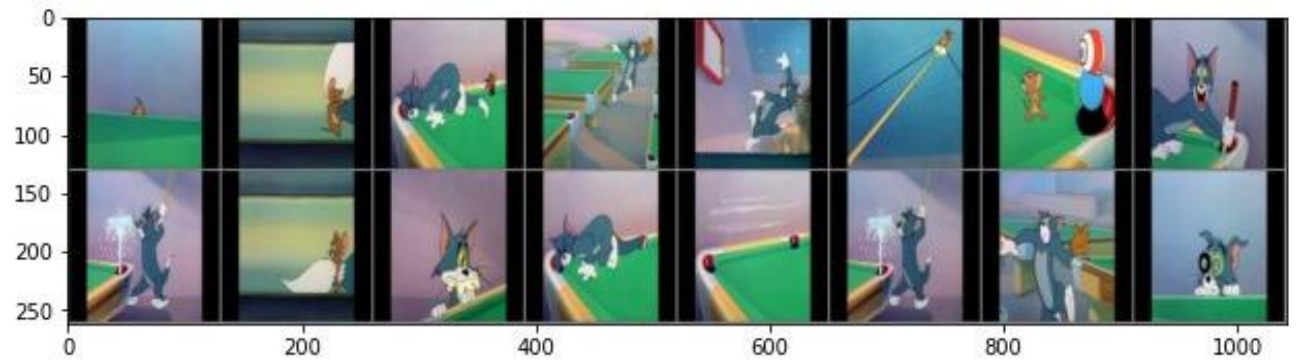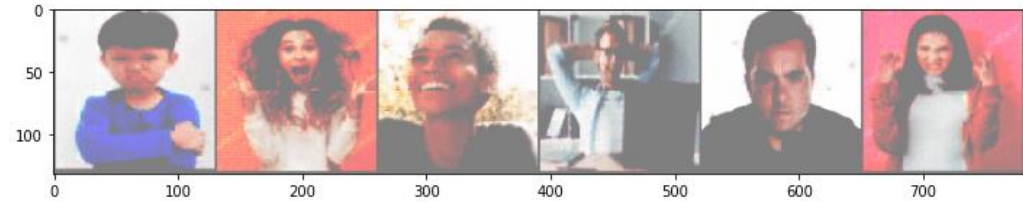# Loss plot of epoch 800-1200

# Training images recovered after training for 1200 epochs

# After training for 1200 epochs, the NN can recover un-seen test images

# Other interesting VAE results

# VAE



$$c_i = exp(\sigma_i) \times e_i + m_i$$

# Encoder

```
[15]:  for batchX, _ in loader:
           break;
       print(batchX.shape)
```

torch.Size([16, 3, 128, 128])

```
(fc1): Linear(in_features=1024, out_features=64,
(fc2): Linear(in_features=1024, out_features=64,
(fc3): Linear(in_features=64, out_features=1024,
```



input → NN Encoder → $m_1$ $m_2$ $m_3$

$\sigma_1$ $\sigma_2$ $\sigma_3$ exp

From a normal distribution $e_1$ $e_2$ $e_3$

$+$ $\times$ → $c_1$ $c_2$ $c_3$

```
16]:  h = model.encoder(batchX.to(device))
      print(h.shape)
```

torch.Size([16, 1024])

```
[17]:  mu=model.fc1(h)
       print(mu.shape)
```

torch.Size([16, 64])

$m_1$ $m_2$ $m_3$

```
[18]:  logvar=model.fc2(h)
       print(logvar.shape)
```

torch.Size([16, 64])

$\sigma_1$ $\sigma_2$ $\sigma_3$

```
[19]:  std = logvar.mul(0.5).exp_()
       print(std.shape)
```

torch.Size([16, 64])

$\sigma_1$ $\sigma_2$ $\sigma_3$ exp

```
[20]:  esp=torch.randn(*mu.size())
       print(esp.shape)
```

torch.Size([16, 64])

$e_1$ $e_2$ $e_3$

```
[21]:  z=mu+std*esp.to(device)
       print(z.shape)
```

torch.Size([16, 64])

$c_1$ $c_2$ $c_3$

# Decoder

```
[22]  z  =  model.fc3(z)
      print(z.shape)
```

torch.Size([16, 1024])

```
[23]  z  =  model.decoder(z)
      print(z.shape)
```

torch.Size([16, 3, 128, 128])

# Loss function



We want $\sigma_i$ close to 0
(variance close to 1)

Minimize

$$\sum_{i=1}^{3} (exp(\sigma_i) - (1 + \sigma_i) + (m_i)^2)$$

L2 regularization

# Loss function

```
[9]:  def loss_fn(recon_x, x, mu, logvar):
        #BCE = F.binary_cross_entropy(recon_x, x, size_average=False).to(device)
        MSE = F.mse_loss(recon_x, x, reduction='sum')
        # see Appendix B from VAE paper:
        # Kingma and Welling. Auto-Encoding Variational Bayes. ICLR, 2014
        # 0.5 * sum(1 + log(sigma^2) - mu^2 - sigma^2)
        KLD = -0.5*torch.mean(1+logvar-mu.pow(2)-logvar.exp()).to(device)
        return MSE+KLD, MSE, KLD
```

Minimize

$$\sum_{i=1}^{3} (exp(\sigma_i) - (1 + \sigma_i) + (m_i)^2)$$

L2 regularization

```
[23]:  tensorY,mu,logvar = model(batchX.to(device))
       print(tensorY.shape)

       torch.Size([16, 3, 128, 128])
```

```
[24]:  loss, mse,kld = loss_fn(tensorY, batchX.to(device), mu, logvar)
       print(loss)

       tensor(627375.3750, device='cuda:0', grad_fn=<AddBackward0>)
```

loss $= MSE(x, \hat{x}) + KL(q(z|x)||P(z))$, why?

loss = $||x - \hat{x}||^2$ + KL[ $N(\mu_x, \sigma_x)$, $N(0, I)$ ] = $||x - d(z)||^2$ + KL[ $N(\mu_x, \sigma_x)$, $N(0, I)$ ]

$$D_{KL}(p||q) = \sum_{i=1}^{N} p(x_i) log(\frac{p(x_i)}{q(x_i)})$$

Minimize

$$\sum_{i=1}^{3} (exp(\sigma_i) - (1 + \sigma_i) + (m_i)^2)$$

19

# The decoder and encoder of VAE model two conditional probabilities

# Gaussian mixture model



**Gaussian Mixture Model**

$$P(x) = \sum_m P(m)P(x|m)$$

How to sample?

$m \sim P(m)$ (multinomial)

m is an integer

$x|m \sim N(\mu^m, \Sigma^m)$

Each x you generate is from a mixture
Distributed representation is better
than cluster.

P(x)

P(m)

1    2    3    4    5    .....

The probability of sampling an output image *x* from latent vector space *z* can be modelled as a Gaussian mixture model



**Gaussian Mixture Model**

$$P(x) = \int_z P(z)P(x|z)dz$$

$P(x|z)$

decoder

$z \rightarrow$ NN $\rightarrow \mu(x)$
$\sigma(x)$

P(x)

Infinite Gaussian

z

We want to train a decoder NN that can maximize the likelihood of observing the training images

## Maximizing Likelihood

$$P(x) = \int_z P(z)P(x|z)dz$$

$$L = \sum_x logP(x)$$

P(z) is normal distribution

$$x|z \sim N\big(\mu(z), \sigma(z)\big)$$

$\mu(z), \sigma(z)$ is going to be estimated

Maximizing the likelihood of the observed x

# Recap: maximize the likelihood of observing the classification of the training data

| Training Data | $x^1$ | $x^2$ | $x^3$ | ... ... | $x^N$ |
|---|---|---|---|---|---|
| | $C_1$ | $C_1$ | $C_2$ | | $C_1$ |

max $\quad L(w,b) = f_{w,b}(x^1)f_{w,b}(x^2)\left(1 - f_{w,b}(x^3)\right)\cdots f_{w,b}(x^N)$

min $\quad -lnL(w,b) = \overset{-}{ln}f_{w,b}(x^1) - lnf_{w,b}(x^2) - ln\left(1 - f_{w,b}(x^3)\right)\cdots$

$\hat{y}^n$: 1 for class 1, 0 for class 2

$$= \sum_n -\left[\hat{y}^n lnf_{w,b}(x^n) + (1 - \hat{y}^n)ln\left(1 - f_{w,b}(x^n)\right)\right]$$

**Cross entropy between two Bernoulli distribution**

Rewrite the maximum likelihood item $\log P(x)$ as the summation of a lower bound $L_b$ and KL divergence

$$logP(x) = \int_z q(z|x)logP(x)dz$$ q(z|x) can be any distribution

$$= \int_z q(z|x)log\left(\frac{P(z,x)}{P(z|x)}\right)dz = \int_z q(z|x)log\left(\frac{P(z,x)}{q(z|x)}\frac{q(z|x)}{P(z|x)}\right)dz$$

$$= \int_z q(z|x)log\left(\frac{P(z,x)}{q(z|x)}\right)dz + \underbrace{\int_z q(z|x)log\left(\frac{q(z|x)}{P(z|x)}\right)dz}_{KL(q(z|x)||P(z|x))}$$ $\geq 0$

$$\geq \int_z q(z|x)log\left(\frac{P(x|z)P(z)}{q(z|x)}\right)dz$$ *lower bound $L_b$*

decoder

$z$ ⟶ NN ⟹ $\mu(x)$ $\searrow$ $\sigma(x)$

$P(x|z)$

encoder

$x$ ⟶ NN' ⟹ $\mu'(x)$ $\searrow$ $z$ $\sigma'(x)$

$q(z|x)$

$$D_{KL}(q||p) = \sum_{i=1}^{N} q(x_i)\log\left(\frac{q(x_i)}{p(x_i)}\right)$$

If we maximum $L_b$ by adjusting $P(x|z)$ and $q(z|x)$ simutaneously, then we can maximum $L_b$ and at the same time minimize the KL distance



$$logP(x) = L_b + KL\big(q(z|x)||P(z|x)\big)$$

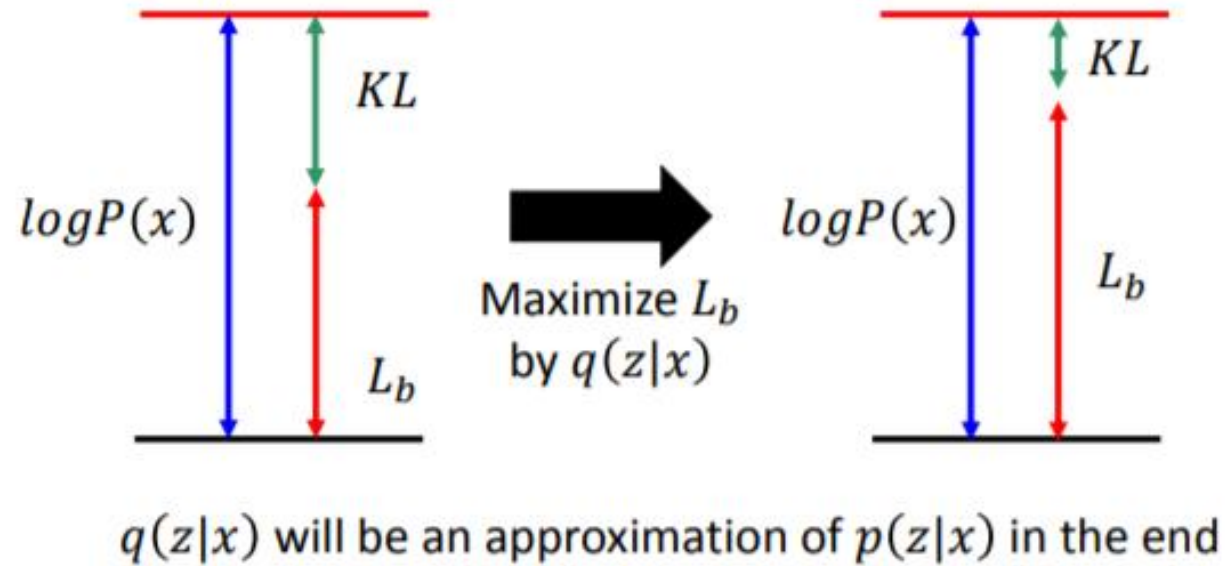$$L_b = \int_z q(z|x) log\left(\frac{P(x|z)P(z)}{q(z|x)}\right) dz$$

Find $P(x|z)$ and $q(z|x)$ maximizing $L_b$

Maximize $L_b$ by $q(z|x)$

$q(z|x)$ will be an approximation of $p(z|x)$ in the end

Rewrite the lower bound $L_b$ as the summation of two terms: $-KL(q(z|x)||P(z))$, and $\int_z q(z|x)\log P(x|z)dz$

max $\quad L_b = \int_z q(z|x)\log\left(\dfrac{P(z,x)}{q(z|x)}\right)dz = \int_z q(z|x)\log\left(\dfrac{P(x|z)P(z)}{q(z|x)}\right)dz$

$= \underbrace{\int_z \boxed{q(z|x)}\log\left(\dfrac{P(z)}{\boxed{q(z|x)}}\right)dz}_{-KL(q(z|x)||P(z))} + \int_z q(z|x)\log P(x|z)dz$

$q(z|x)$

encoder

$x \Rightarrow$ [NN'] $\nearrow \mu'(x) \searrow$ $z$
$\searrow \sigma'(x) \nearrow$

$P(x|z)$

decoder

$z \Rightarrow$ [NN] $\nearrow \mu(x)$
$\searrow \sigma(x)$

Reference: 李弘毅 ML Lecture 18  https://youtu.be/8zomhgKrsmQ

max. $L_b$ can be done by min. $KL(q(z|x)||P(z))$ and max. $\int_z q(z|x)logP(x|z)dz$. That is why loss = KLD + MSE (x, $\hat{x}$)

Minimizing $KL(q(z|x)||P(z))$

Minimize

$$\sum_{i=1}^{3}(exp(\sigma_i) - (1 + \sigma_i) + (m_i)^2)$$

$x \rightarrow$ NN' $\nearrow \mu'(x)$ $\searrow \sigma'(x)$

(Refer to the Appendix B of the original VAE paper)

Maximizing

$$\int_z q(z|x)logP(x|z)dz = E_{q(z|x)}[logP(x|z)]$$

By min. $MSE(x, \hat{x})$

close

$x \rightarrow$ NN' $\nearrow \mu'(x)$ $\searrow \sigma'(x)$ $\rightarrow z \rightarrow$ NN $\nearrow \mu(x) \leftrightarrow x$ $\searrow \sigma(x)$

# HW6 (2)

- Train an VAE to learn a compact representation (try latent vector of size 20, 30, 50) of your facial expression. Test with 10 happy and 10 angry faces.

- Show the recovered image.

- Send the latent vectors to $t$-SNE to see whether they form clusters.