# 1   Paper Introduction

This homework explores Latent Diffusion Models as presented in the paper High Resolution Image Synthesis with Latent Diffusion. Stable diffusion is a text conditional latent diffusion model, created by researchers and engineers from CompVis, Stability AI, and LAION. Diffusion models are incredibly expensive to the train so with the use of latent space, we can greatly reduce the computational resource requirements for training, and also result very competitive performance on tasks such as inpainting, text-to-image synthesis, super-resolution and more. We encourage you to read over the paper on your own as it is quite fascinating!

## 1.1   Auto-encoders and Downsampling

Latent diffusion models use pre-trained auto-encoders to produce a lower dimensional space representation of the input. Auto-encoders downsample the input images by a factor of $f = 2^m, m \,\epsilon\, \mathbb{R}$. The paper experiments with 6 different sampling factors $f = 1, 2, 4, 8, 16, 32$ abbrieviated as LDM-f. LDM-1 represents pixel-based diffusion models with no downsampling. The paper tested the performance on these 6 models and measured performance with two different scores: FID (Frechet Inception Distance) which compares the distributions of images generated with the model and real images and IS (Inception Score) which quantifies the quality of images generated by the model.

(a) Given two graphs that show the FID scores and IS scores against training progress, we can see that three of the LDM-f models perform well and three have weaker performance. Identify the top three performing LDM-f models and the bottom three performing LDM-f models. Provide a brief explanation for your answer.
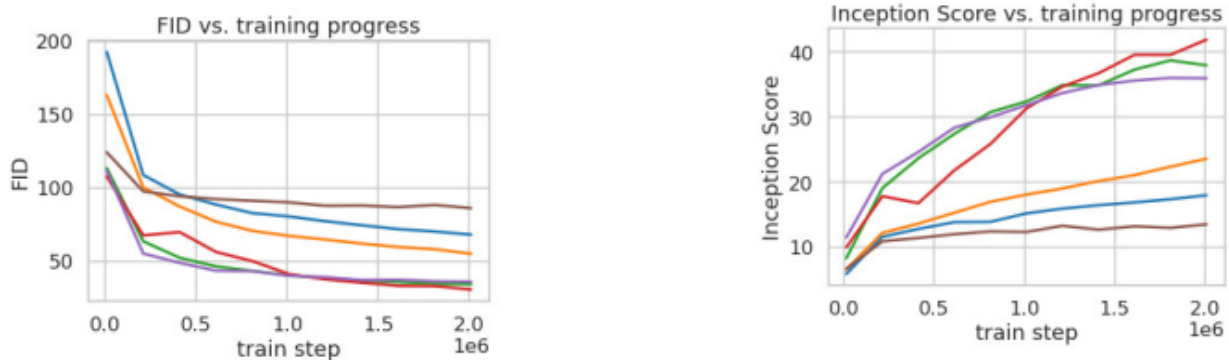


Abbildung 1: IS Scores vs Training Progress

## 1.2   Conditional LDM Loss

We can think of diffusion models as a sequence of denoising autoencoders $\epsilon_\theta(x_t, t)$; $t = 1...T$ where $x_t$ is a noisy version of the input and $t$ is sampled uniformly. We get the following loss:

$$L_{LDM} = \mathbb{E}_{\varepsilon(x), \varepsilon \sim \mathcal{N}(0,1), t}[\|\varepsilon - \varepsilon_\theta(x_t, t)\|_2^2]$$

For a conditional LDM, we use a conditional denoising autoencoder which takes in the latent representation, $z_t$ and a domain specific encoder $\tau_\theta$ that projects a text input $y$ into an intermediate representation.

(b) Write the loss function for a conditional LDM.

# 2   Image Synthesis with Diffusion Model

Please follow the instructions in this notebook. You will implement a text conditioned denoising unet and the sampling algorithm. Then you'll see the diffusion process on a single image and try to train your own diffusion mode. Once you finished with the notebook,

- download submission.zip and submit it to Gradescope.

- Answer the following questions in your submission of the written assignment:

(a) Why might a U-Net by a good choice for the model backbone? List two reasons.

(b) Screenshot your visualization for get noisy image and include it in your submission of the written assignment. And describe the picture you observed briefly. What kind of process is this?



Abbildung 2: Visualization for Sampling

(c) Screenshot one of your visualizations about the sampled images and include it in your submission of the written assignment. Then answer the following question. How does the model perform and does it meet your expectations? If not, what do you think are the directions for improvement?
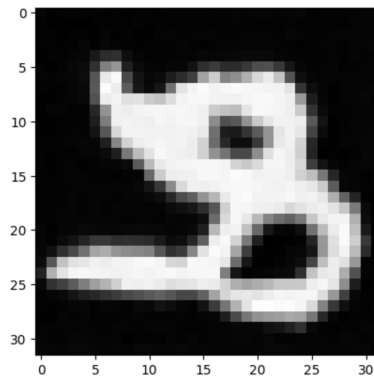


Abbildung 3: Visualization for Model Performance

# 3 Understanding the VAE Loss

Latent diffusion models use an autoencoder model similar to a VAE to encode images from pixel space to some compressed latent space, perform the diffusion process in this latent space, and finally decode the denoised latent back into pixel space.

To have our model learn the correct latent space, we would like the estimated posterior by the encoder $q_\phi(z|x_i)$ to be very close to the real posterior $p_\theta(z|x_i)$. We can quantify this using the Kullback-Leibler (KL) divergence to get a measure of the difference between the two distributions, defined as

$$D_{KL}(P||Q) = \mathbb{E}_{x \sim P(x)}[log\frac{P(x)}{Q(x)}]$$

However, $p_\theta(z|x_i)$ is intractable in practice. Another goal of VAE-like models is to maximize the probability of generating real samples, or $p_\theta(x)$. During training, rather than directly maximizing this quantity, we maximize a lower-bound called the evidence lower bound (ELBO):

$$log(p_\theta(x_i)) \geq \mathcal{L}(x_i, \theta, \phi) = \mathbb{E}_{z \sim q_\phi(z|x_i)}[log(p_\theta(x_i|z)] - D_{KL}[q_\phi(z|x_i)||p_\theta(z)]$$

At a high-level, this loss can be seen as doing two things:

- **Reconstruction loss**: The first term corresponds to the likelihood of generating images from the true data distribution given a sampled latent

- **Matching Prior**: The second term acts as a regularizer to minimize the difference between the estimated posterior and the prior (generally a standard normal).

Note: In the latent diffusion paper, the autoencoder loss is actually comprised of 3 terms: a patch-based discriminator loss for reconstruction, the standard KL-divergence term between $q_\phi(z|x)$ and a prior $p(z)$, and a regularizing term $L_{reg}$ for the latent z to be zero centered and with small variance. For simplicity, however, we will just be analyzing the traditional ELBO loss here.

(a) **Show** that maximizing the ELBO loss also minimizes $D_{KL}(q_\phi(z|x_i)||p_\theta(z|x_i))$, the KL-divergence of the estimated posterior and the real posterior.

*Hint: Start with the KL-divergence formula and manipulate it to get ELBO loss plus some other term.*