



# **Hands-on Session on Image Generative Models with Dall.E Mini**



AI model drawing images from any prompt!

the Eiffel tower landing on the moon

DRAW



# Dall.E Mini — Text to Image

[Online Live Version of Dall.E Mini](#)



## Part 1: Building Blocks of Dall.E Mini

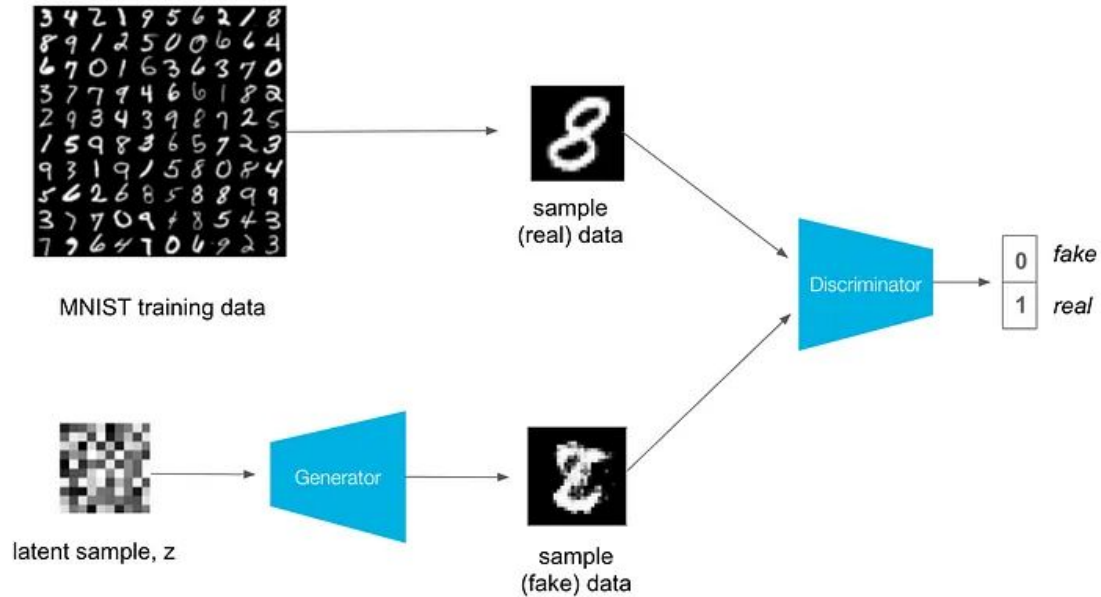
- **BERT-based Encoder-Decoder:** Encodes captions as embedding vectors
- **VQ-GAN:** Decodes caption embeddings into Images
- **CLIP:** Evaluates caption-image relevance



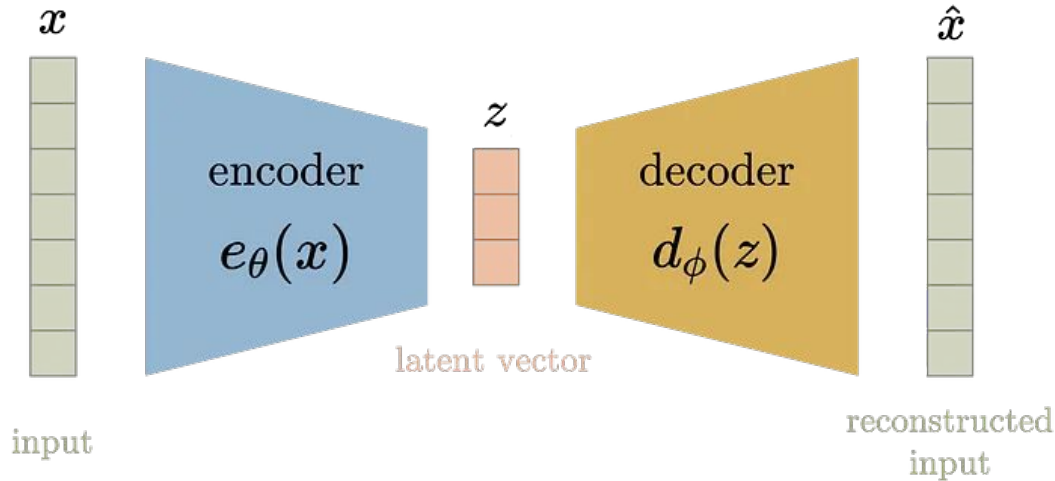
## Part 2: Generative Adversarial Networks (GANs)

- Dall.E Mini uses a variant of GANs called VQ-GANs.
- The evolution of VQ-GANs,
  - Simple GAN
  - Autoencoders (AEs)
  - Variational Autoencoders (VAEs)
  - Vector Quantized Autoencoders (VQ-AEs)
  - Vector Quantized GANs (VQ-GANs)

# Simple GAN



## Autoencoder (AE)



- The latent space is discontinuous and has significant “gaps”.

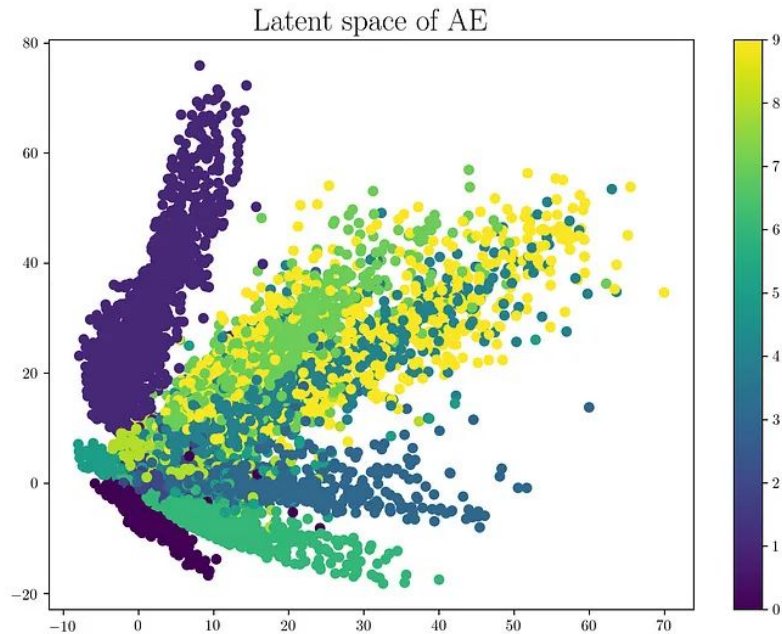
---

$$loss = \|x - \hat{x}\|_2 = \|x - d_{\phi}(z)\|_2 = \|x - d_{\phi}(e_{\theta}(x))\|_2$$

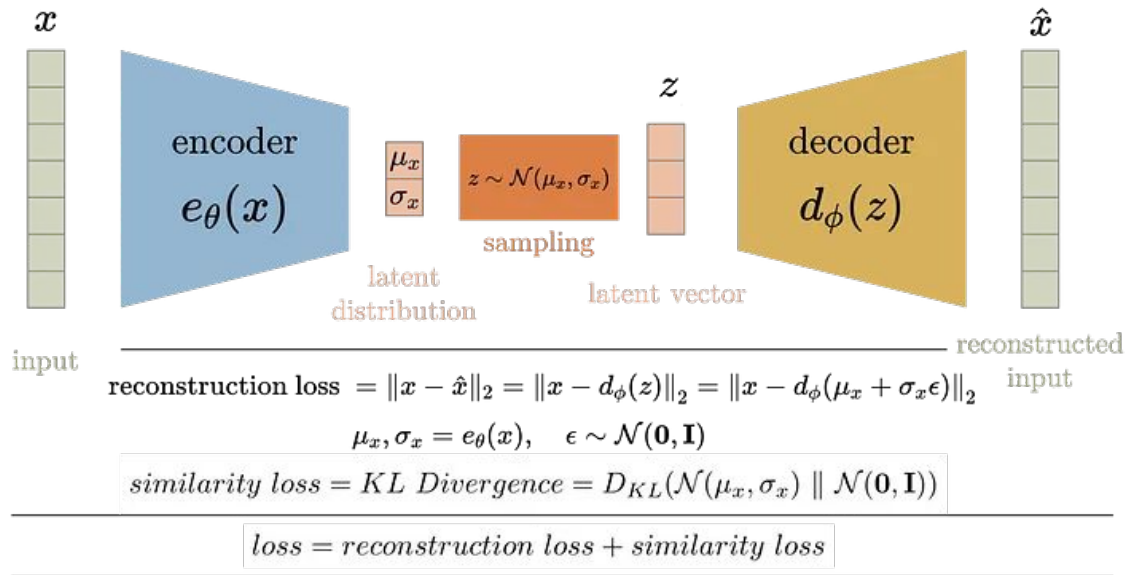
---



# Autoencoder (AE)



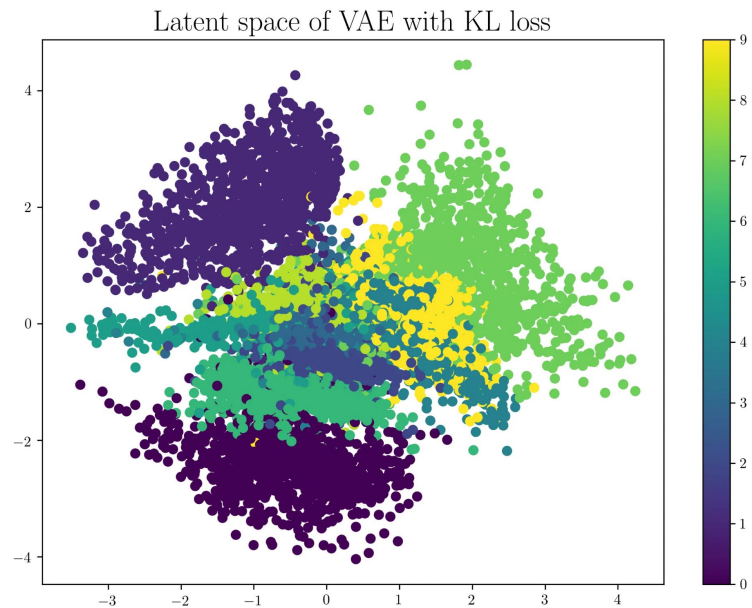
# Variational Autoencoder (VAE)



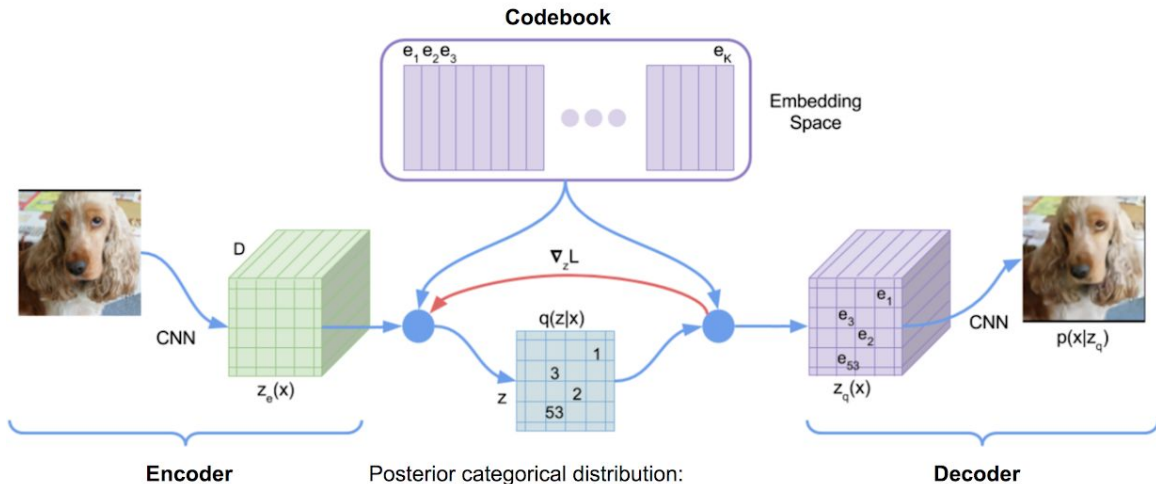
- The latent space is more cohesive – resembles the unit norm.
- Overlapping regions produce “morphed” images.



# Variational Autoencoder (VAE)



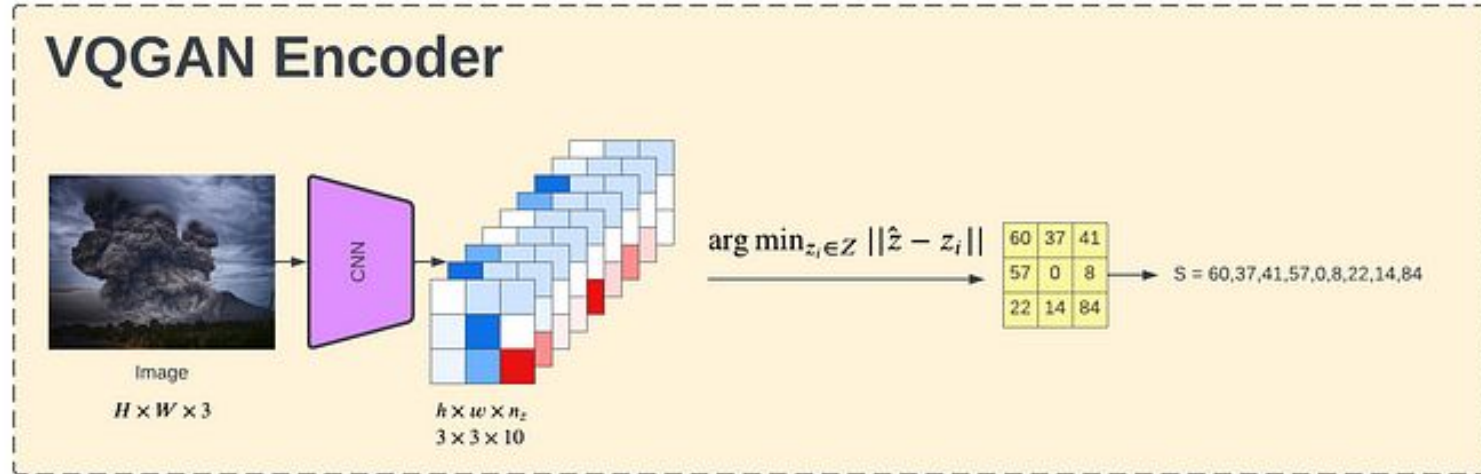
# Vector-Quantized Variational Autoencoder (VQ-VAE)



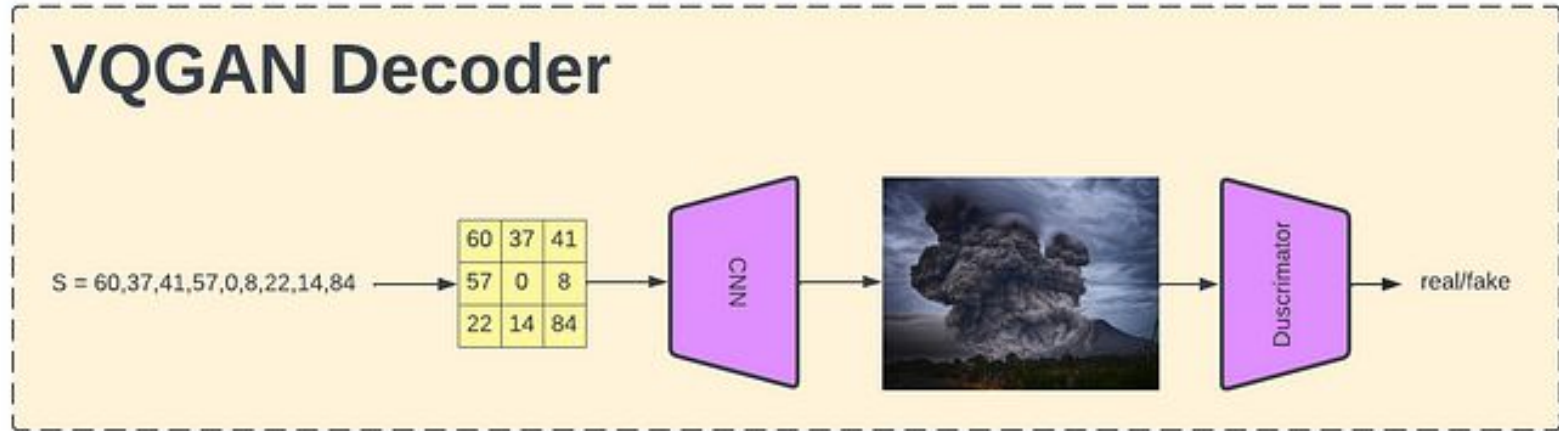
- The latent space is discrete.
- No “morphed” outputs.
- Latent space has same dimensions as codebook.

$$q(\mathbf{z} = \mathbf{e}_k | \mathbf{x}) = \begin{cases} 1 & \text{if } k = \arg \min_i \|\mathbf{z}_e(\mathbf{x}) - \mathbf{e}_i\|_2 \\ 0 & \text{otherwise.} \end{cases}$$

# Vector-Quantized Variational GAN (VQ-GAN)



## Vector-Quantized Variational GAN (VQ-GAN)

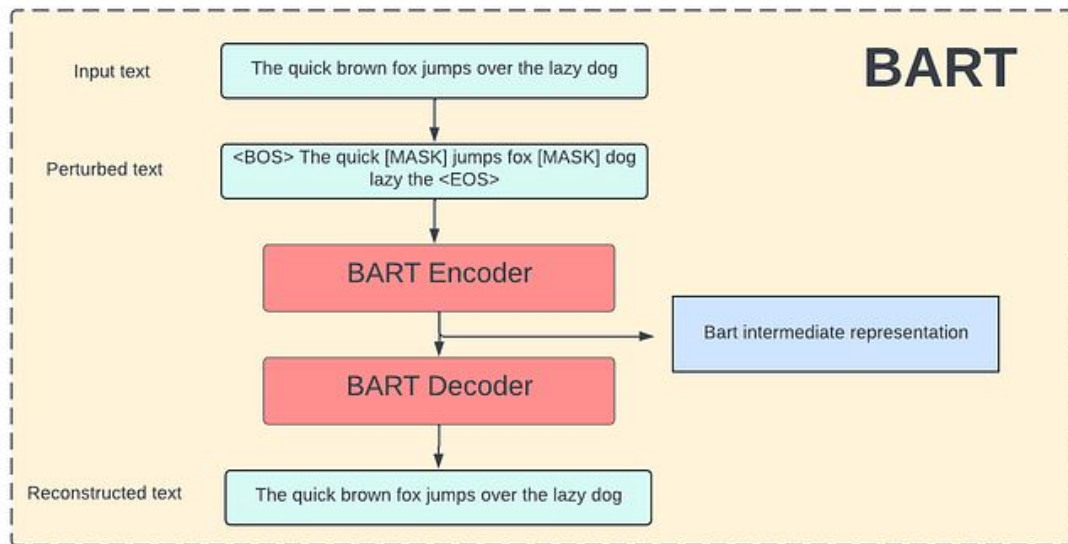




## Part 3: BART Encoder-Decoder

- A BART model is pre-trained to “clean” text captions.
- For Dall.E Mini, the BART model **translates captions into the codebook vocabulary**.
- The codebook of VQ-GAN, in effect, maps text embeddings to image embeddings.

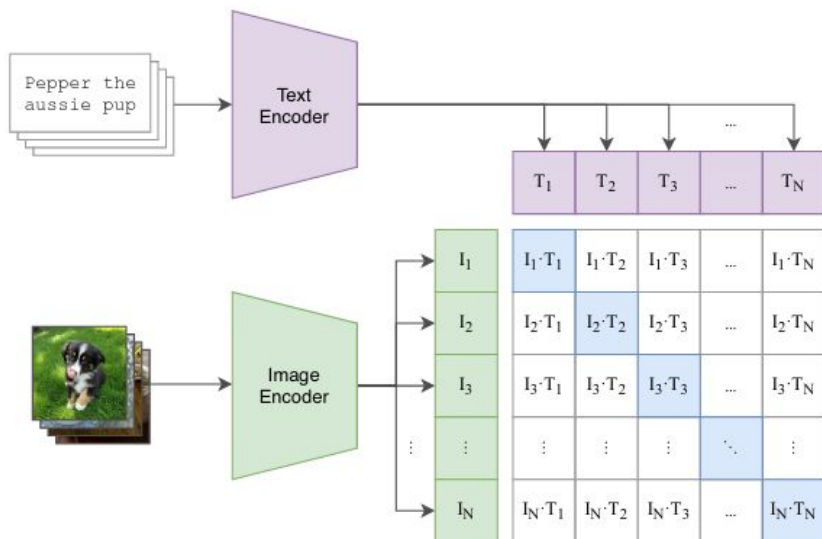
## What BART does for Dall.E.



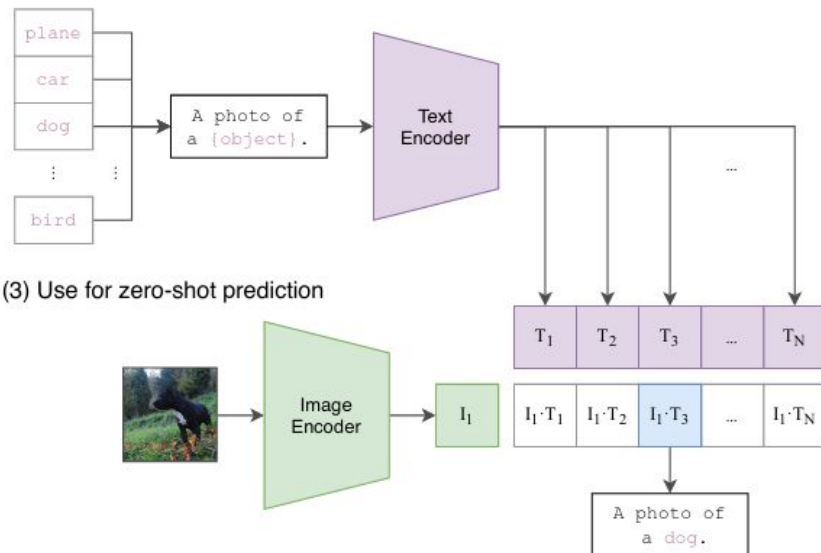
- Translates captions to codebook vocabulary.

## Part 4: CLIP to Rank Images by Relevance

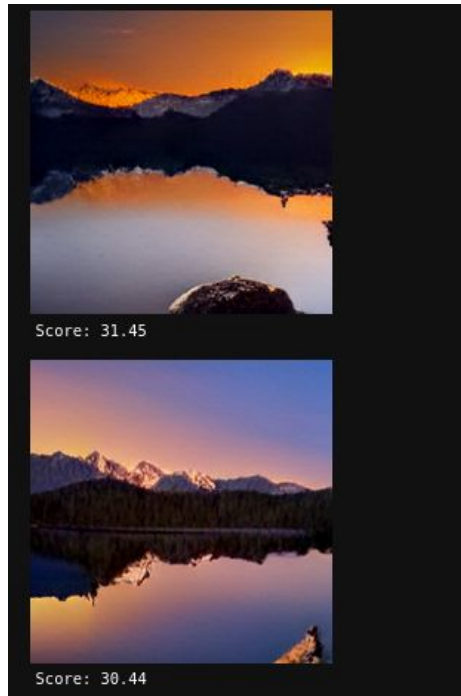
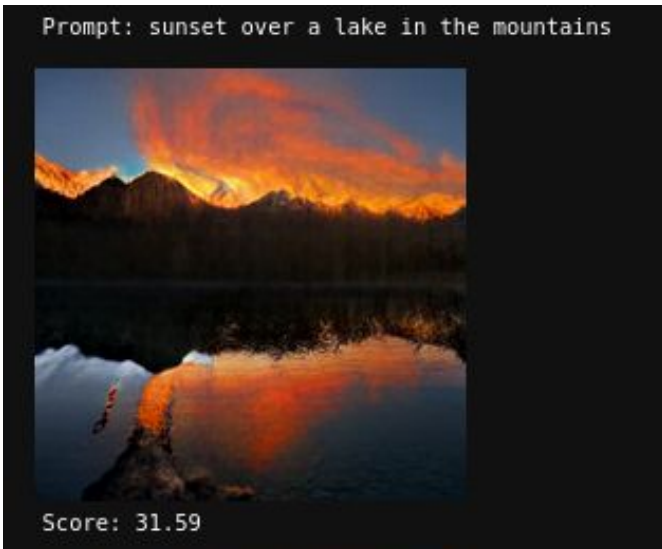
(1) Contrastive pre-training



(2) Create dataset classifier from label text

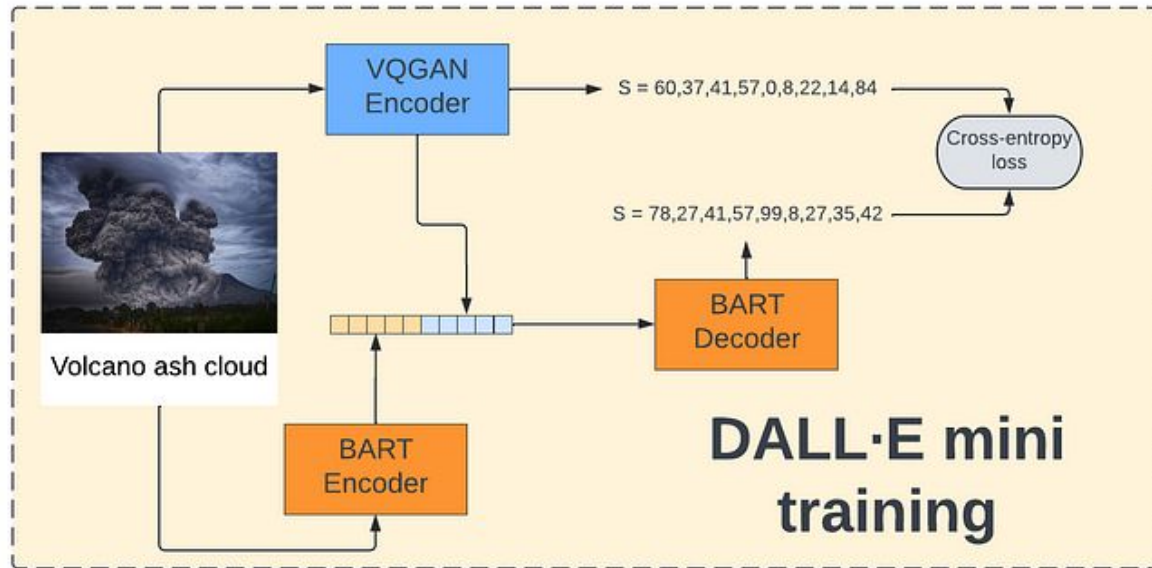


## Part 4: CLIP to Rank Images by Relevance

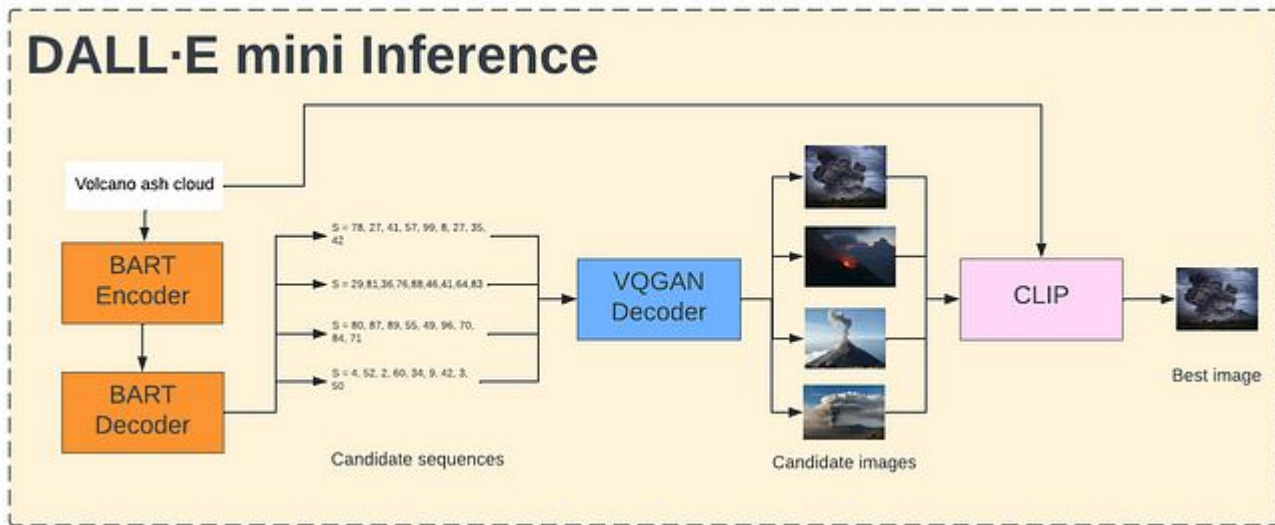




## Part 5: Putting the pieces together.



# The Dall.E Mini Text-to-Image Pipeline.



# Examples of Generated Images

TEXT PROMPT

a stained glass window with an image of a blue strawberry

AI-GENERATED  
IMAGES



# Examples of Generated Images



# Examples of Generated Images

TEXT PROMPT a photo of the food of china

AI-GENERATED  
IMAGES

