

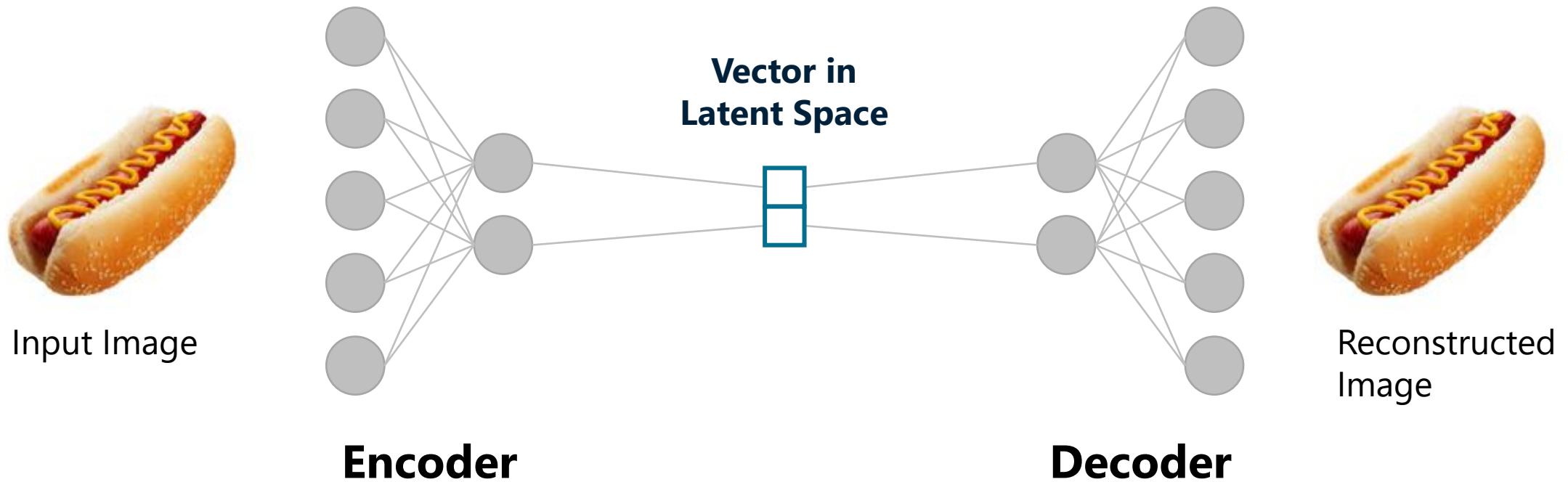
# IMAGE GENERATION

Jeff Prosise

[jeff.prosise@atmosera.com](mailto:jeff.prosise@atmosera.com)

# Autoencoders

- Use encoders to reduce inputs to lower-dimensional "latent space" and decoders to reconstruct the original inputs from latent space



# Implementing an Autoencoder

```
encoder = Sequential()  
encoder.add(Flatten()) # Reshape 28x28 images to 1D  
encoder.add(Dense(128, activation='relu'))  
encoder.add(Dense(32, activation='relu'))  
  
decoder = Sequential()  
decoder.add(Dense(128, activation='relu'))  
decoder.add(Dense(28 * 28, activation='relu'))  
decoder.add(Reshape((28, 28))) # Reshape 1D array into 28x28 image  
  
model = Sequential([encoder, decoder])  
model.compile(loss='mse', optimizer='adam')  
model.fit(x_train, x_train, epochs=10, validation_data=(x_test, x_test))
```

# Implementing a Convolutional Autoencoder

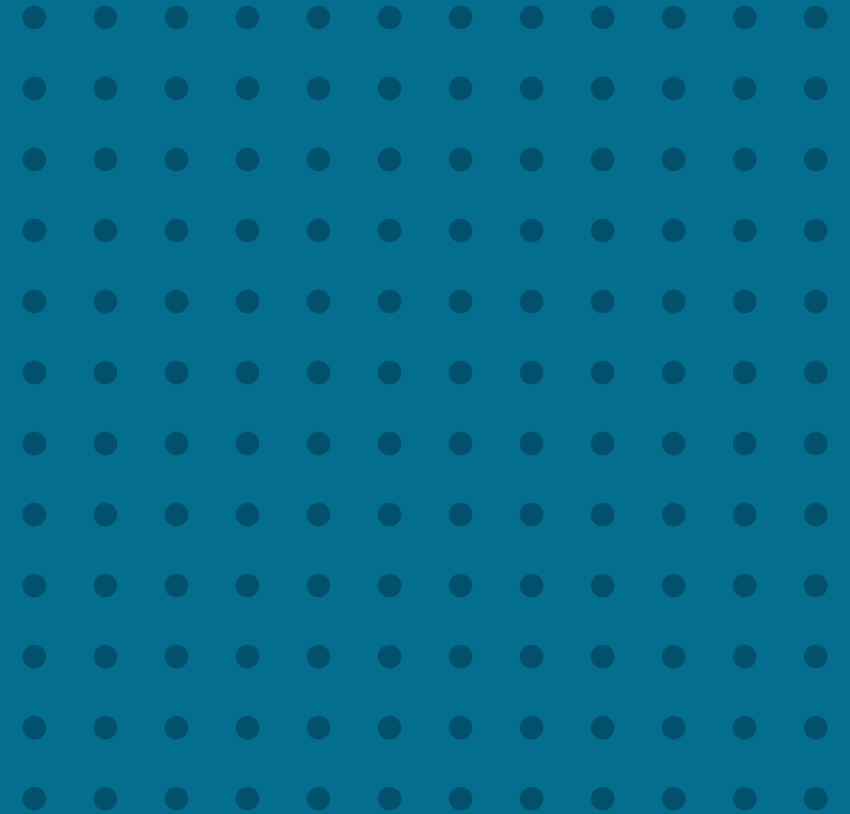
```
encoder = Sequential()
encoder.add(Conv2D(16, (3, 3), activation='relu', strides=2, padding='same',
                  input_shape=(28, 28, 1)))
encoder.add(Conv2D(32, (3, 3), activation='relu', strides=2, padding='same'))
encoder.add(GlobalAveragePooling2D())

decoder = Sequential()
decoder.add(Dense(7 * 7 * 16), activation='relu')
decoder.add(Reshape((7, 7, 16)))
decoder.add(Conv2DTranspose(32, (3, 3), strides=2, activation='relu', padding='same'))
decoder.add(Conv2DTranspose(1, (3, 3), strides=2, activation='relu', padding='same'))
decoder.add(Reshape((28, 28)))

model = Sequential([encoder, decoder])
model.compile(loss='mse', optimizer='adam')
model.fit(x_train, x_train, epochs=10, validation_data=(x_test, x_test))
```

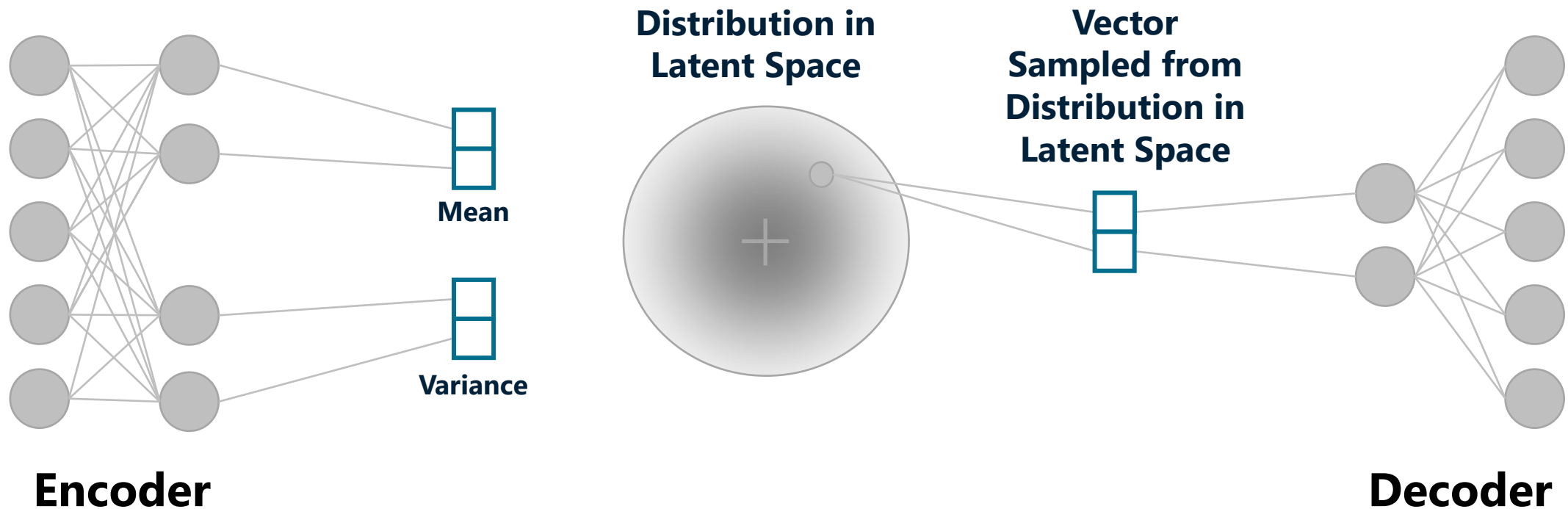
# Demo

Autoencoders



# Variational Autoencoders (VAEs)

- Provide continuity in latent space by introducing randomness
- Encoder generates probabilistic distributions rather than points



# Implementing a Variational Encoder

```
encoder_input = Input(shape=(28, 28))
x = Flatten()(encoder_input)
x = Dense(128, activation='relu')(x)
x = Dense(32, activation='relu')(x)
z_mean = Dense(8)(x) # Outputs vector representing mean in normal distribution
z_log_var = Dense(8)(x) # Outputs vector representing variance in normal distribution

encoder_output = Sampling()([z_mean, z_log_var]) # Custom layer with two inputs that
                                                # outputs vector sampled from normal
                                                # distribution defined by the inputs

encoder = Model(encoder_input, encoder_output)
```

# Computing Loss

- Loss is the sum of reconstruction loss (MSE) and latent loss (Kullback-Liebler divergence between target distribution and actual distribution)

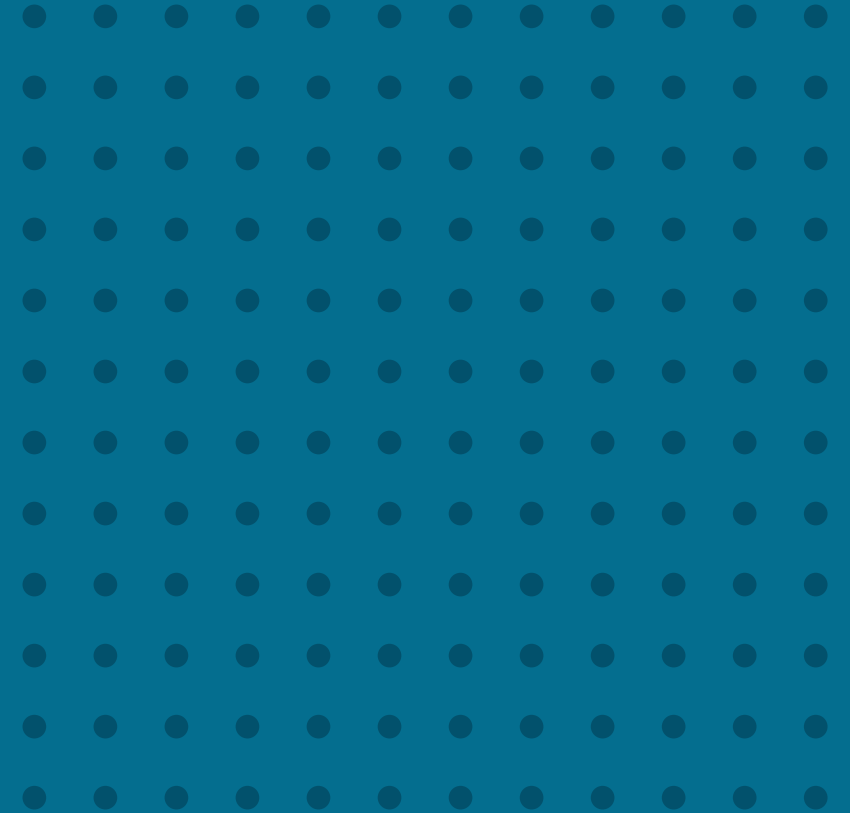
$$\mathcal{L} = -\frac{1}{2} \sum_{i=1}^n [1 + \gamma_i - \exp(\gamma_i) - \mu_i^2]$$

```
model = Model(encoder_input, decoder(encoder(encoder_input)))
model.compile(loss='mse', optimizer='adam')
# Add a KL loss function to the model's loss computations
z_loss = -0.5 * tf.reduce_sum(1 + z_log_var - tf.exp(z_log_var) -
                             tf.square(z_mean), axis=-1)
model.add_loss(tf.reduce_mean(z_loss) / (28 * 28))
model.fit(x_train, x_train, epochs=10, validation_data=(x_test, x_test))
```

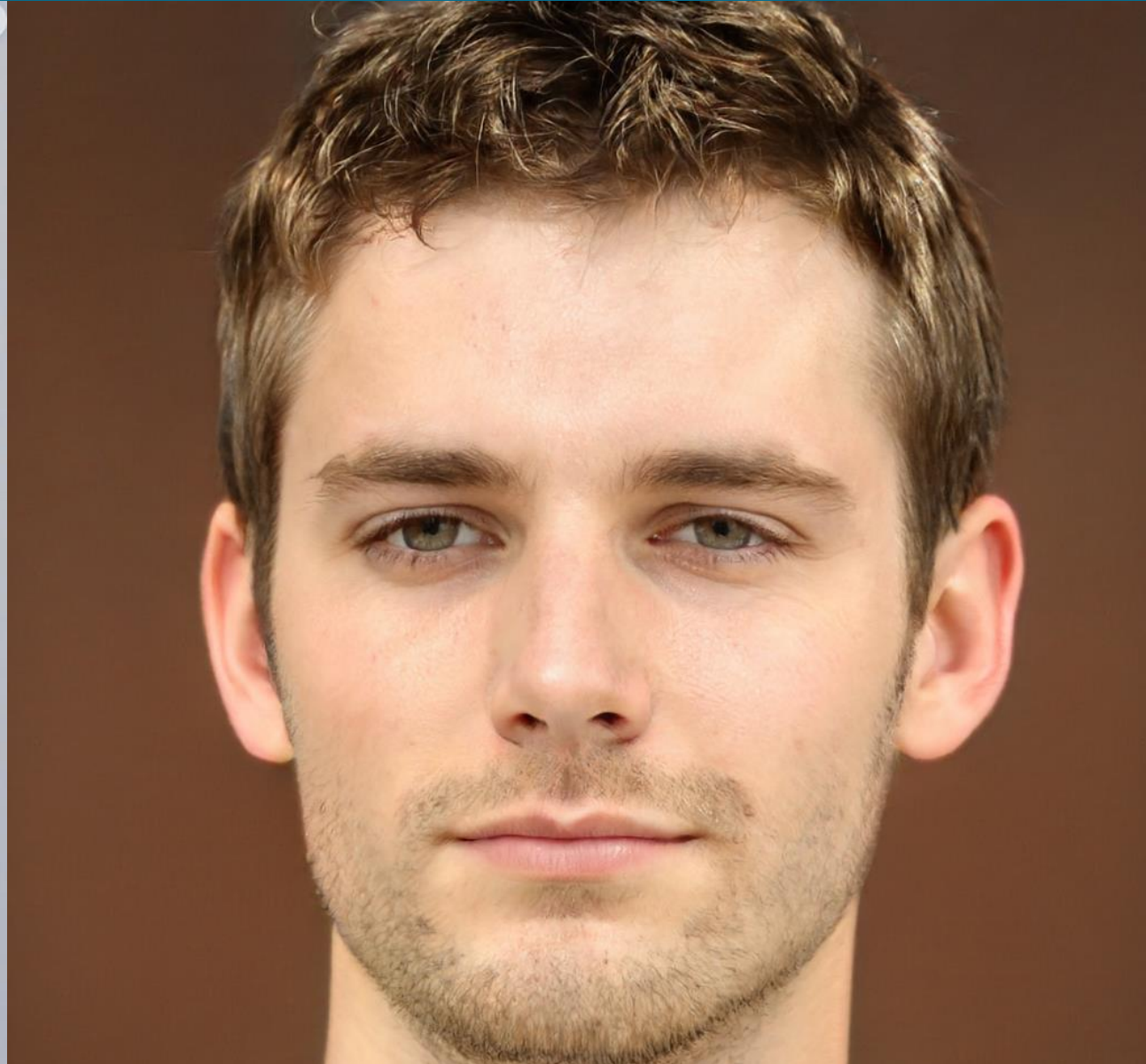


# Demo

Variational Autoencoders

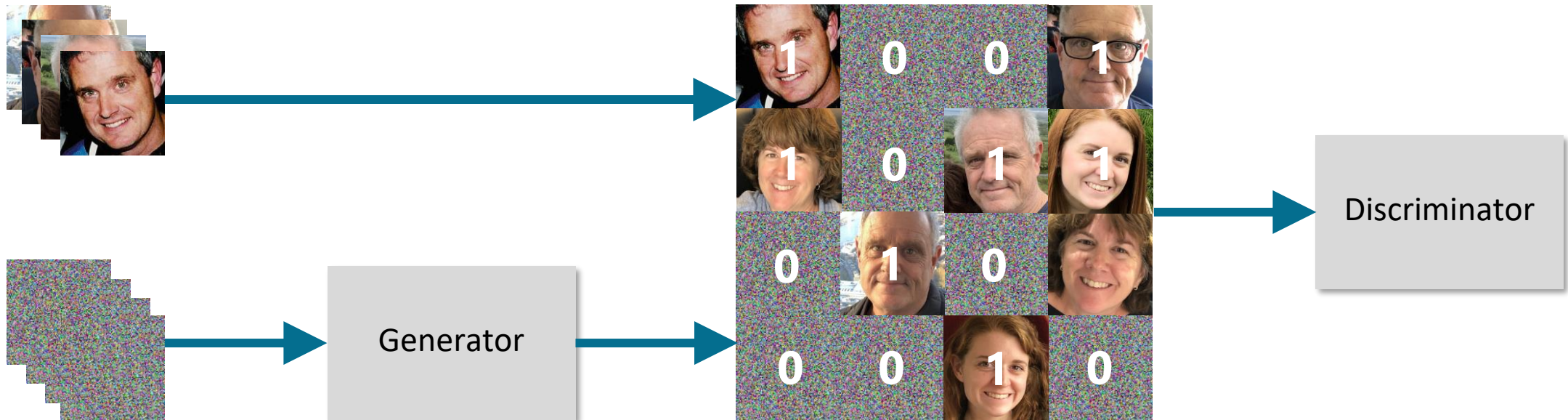


[thispersondoesnotexist.com](http://thispersondoesnotexist.com)



# Generative Adversarial Networks (GANs)

- Use one model (the *discriminator*) to train another model (the *generator*) to produce images from random inputs
- After training, use the generator to produce realistic images



# Building and Training GANs

- Architect model to avoid mode collapse and instability
  - Use strided convolutions (strides=2) rather than pooling and upsampling layers in both the generator and the discriminator
  - Use batch normalization after all layers except for the generator's output layer and the discriminator's input layer
  - Use ReLU activation in the generator, except for the last layer, which should use **tanh** activation instead (requires preconditioning of training data)
  - Use leaky ReLU activation in the discriminator
- Implement custom training loop to train generator and discriminator separately in each batch of each epoch

# Implementing a Discriminator

```
discriminator = Sequential()  
discriminator.add(Conv2D(16, (3, 3), activation=LeakyReLU(0.2), strides=2,  
                        padding='same'))  
discriminator.add(BatchNormalization()) # Or Dropout  
discriminator.add(Conv2D(32, (3, 3), activation=LeakyReLU(0.2), strides=2,  
                        padding='same'))  
discriminator.add(BatchNormalization()) # Or Dropout  
discriminator.add(Flatten()) # Or GlobalAveragePooling2D  
discriminator.add(Dense(1, activation='sigmoid'))  
discriminator.compile(loss='binary_crossentropy', optimizer='adam')  
discriminator.trainable = False
```



# Implementing a Generator

```
generator = Sequential()
generator.add(Dense(7 * 7 * 16))
generator.add(Reshape([7, 7, 16]))
generator.add(BatchNormalization())
generator.add(Conv2DTranspose(32, (3, 3), strides=2, activation='relu', padding='same'))
generator.add(BatchNormalization())
generator.add(Conv2DTranspose(16, (3, 3), strides=2, activation='tanh', padding='same'))

# Form a GAN from the generator and the discriminator
model = Sequential([generator, discriminator])
model.compile(loss='binary_crossentropy', optimizer='adam')

# Don't call fit(); model requires a custom training loop
```

# Training the Model

```
for epoch in range(epochs):  
    np.random.shuffle(x_train) # Shuffle images at the start of each epoch  
  
    for step in steps_per_epoch:  
        # Train the discriminator to differentiate between real and generated images  
        # TODO: Fetch a batch of real images and combine with a batch of fake images (x)  
        # TODO: Generate labels for images where 0 == fake and 1 == real (y)  
        discriminator.train_on_batch(x, y)  
  
        # Train the model (trains the generator but not the discriminator)  
        x = np.random.normal(size=(batch_size, input_size))  
        y = np.array([1] * batch_size)  
        model.train_on_batch(x, y)
```

# Demo

Generative Adversarial Networks

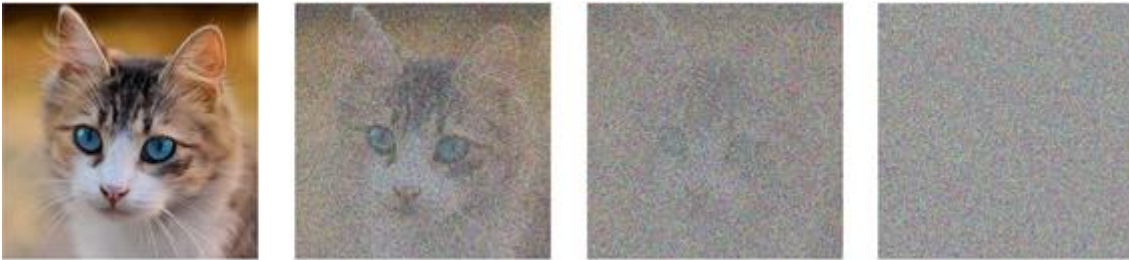




# Diffusion Models

- Trained to turn random noise into images using diffusion process
- Commercial examples include Stable Diffusion, ImageGen, and DALL·E
- Slower at inference time because decoding (like training) is progressive

## Training



During training, images have **random noise** added in stages. The model is trained to **remove the noise** one step at a time, reversing the progression.

## Inference (Generation)



At inference time, an image **filled with random pixel values** is input to the model. The model repeatedly denoises the image to produce a final image.

# Text-Guided Diffusion



## Prompt

A hyperrealistic photograph of ancient Tokyo/London/Paris architectural ruins in a flooded apocalypse landscape of dead skyscrapers, lens flares, cinematic, hdri, matte painting, concept art, celestial, soft render, highly detailed, cgsociety, octane render, trending on artstation, architectural

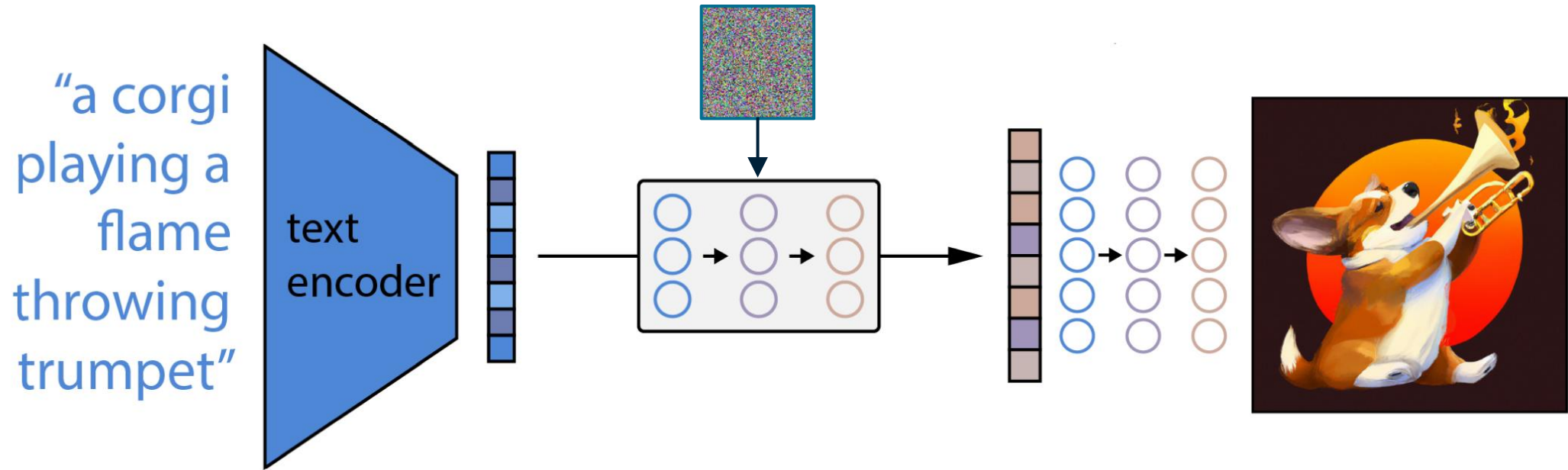


# DALL·E 2

- Diffusion model from OpenAI featuring 3.5 billion parameters
- Trained on 650 million text-image pairs scraped from the Internet
- Available by REST API
  - Requires an OpenAI account
- Supports image generation and modification, image variations, inpainting, and outpainting



# How DALL·E 2 Works



Contrastive Language-Image Pretraining (CLIP) model encodes the prompt, **generating a text embedding** in latent space

"Prior" model generates an **image embedding** from the text embedding and random noise

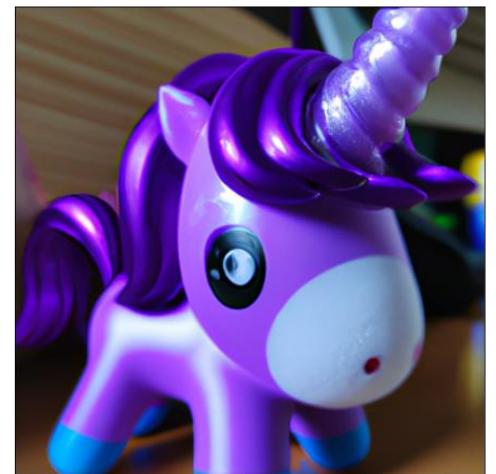
Decoder uses **reverse diffusion** to **generate a 64x64 image** from the image embedding. CNNs **upsample the image** to 256x256, 512x512, and 1,024x1,024 to produce the final image.

# Generating Images with DALL·E 2

```
import openai, base64, PIL, io
openai.api_key = 'OPENAI_API_KEY'

response = openai.Image.create(
    prompt='Photo of a purple unicorn',
    size='512x512',
    n=1,
    response_format='b64_json'
)

image_data = response['data'][0]['b64_json']
image = io.BytesIO(base64.b64decode(image_data))
```



# Creating Variations of Existing Images

```
response = openai.Image.create_variation(  
    image=open('PATH_TO_IMAGE', 'rb'),  
    size='512x512', n=1, response_format='b64_json'  
)
```

Original image



Variation created  
by DALL·E 2



# Inpainting

```
response = openai.Image.create_edit(  
    image=open('PATH_TO_IMAGE', 'rb'), # Path to original image  
    mask=open('PATH_TO_IMAGE_MASK', 'rb'), # Path to same image with transparent pixels  
    prompt='Photograph of two people standing on a cliff overlooking the beach',  
    n=1, size='512x512', response_format='b64_json'  
)
```

Original image  
with fence in  
background

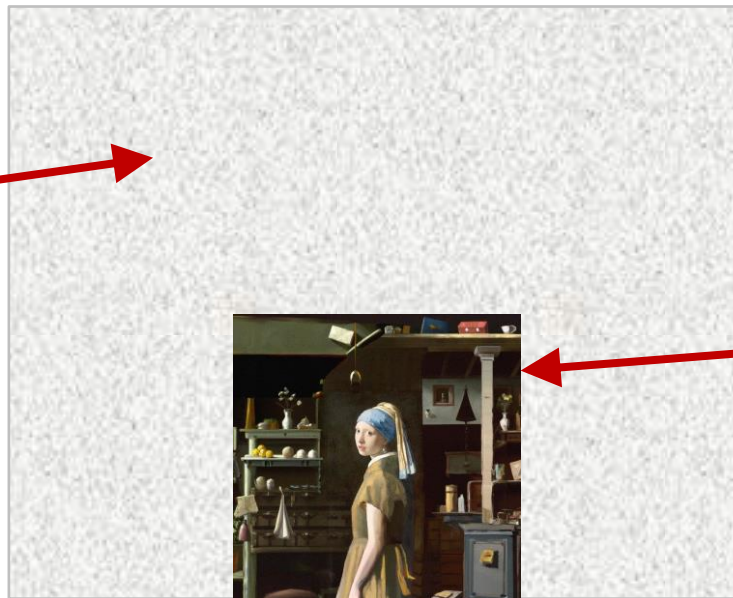


Image mask with  
transparent pixels  
denoting regions  
to be inpainted

# Outpainting

```
response = openai.Image.create_edit(  
    image=open('PATH_TO_IMAGE', 'rb'),  
    mask=open('PATH_TO_IMAGE', 'rb'), # Same image  
    prompt='Painting of a girl standing in a kitchen',  
    n=1, size='512x512', response_format='b64_json'  
)
```

Transparent pixels  
identifying region  
to be outpainted



Region to be  
expanded via  
outpainting



# Demo

DALL·E 2

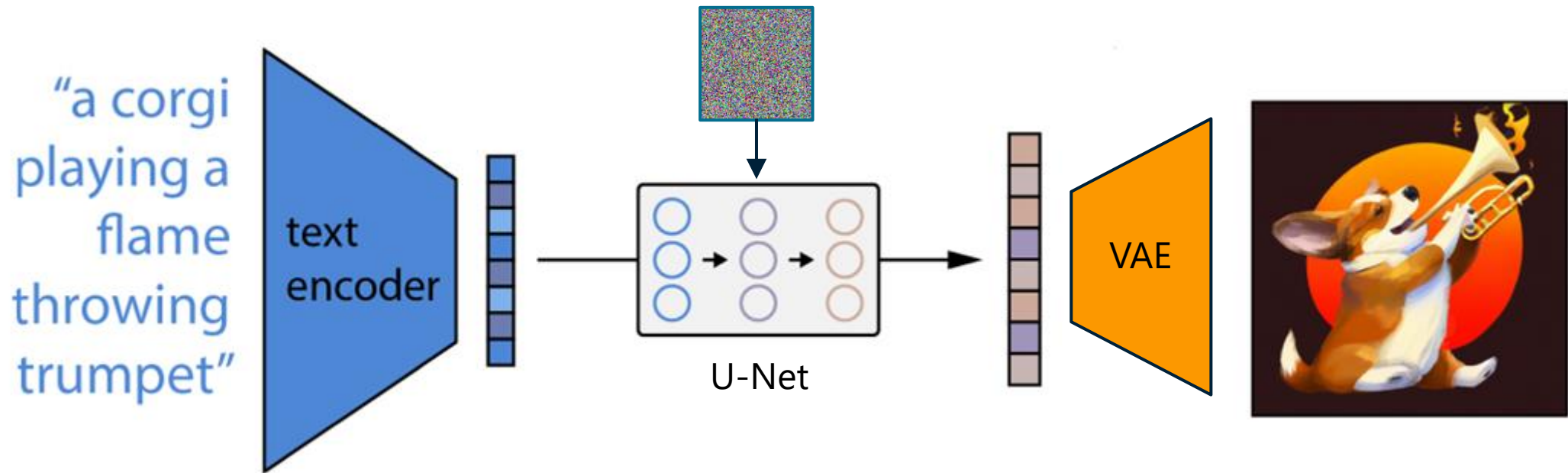


# Stable Diffusion

- Text-guided latent diffusion model published by Stability AI
  - Faster and less memory hungry than models that work in pixel space
  - Supports image generation (any size), inpainting, and outpainting
- Open-sourced (code and weights) in August 2022



# How Stable Diffusion Works



Contrastive Language-Image Pretraining (CLIP) model based on ClipText (V1) or OpenCLIP (V2) encodes the prompt, **generating a text embedding** in latent space

**U-Net** reverse diffusion model generates an **image embedding** from the text embedding and random noise. Diffusion is performed in **latent space** that is **48 times smaller** than pixel space.

Variational autoencoder **generates a 64x64 image** from the image embedding. Image is **upsampled** to 512x512 (V1) or 768x768 (V2) to produce the final image.

# Generating Images with Stable Diffusion V1

```
from keras_cv.models import StableDiffusion
```

```
model = StableDiffusion(img_width=900, img_height=500)
```

```
image = model.text_to_image('Photograph of a city street', batch_size=1)[0]
```

# Generating Images with Stable Diffusion V2

```
from keras_cv.models import StableDiffusionV2

model = StableDiffusionV2(img_width=900, img_height=500)
image = model.text_to_image('Photograph of a city street', batch_size=1)[0]
```

# Inpainting with Stable Diffusion V2

```
model = StableDiffusionV2(  
    img_width=896, # Must be multiple of 32  
    img_height=512 # Must be multiple of 32  
)  
  
image = model.inpaint(  
    'Photo of a city street',  
    negative_prompt='Remove person from the photo',  
    image=photo, # 3-channel RGB image to be inpainted  
    mask=mask,   # Mask of 1s and 0s with 0s denoting regions(s) to be inpainted  
    batch_size=1  
)[0]
```

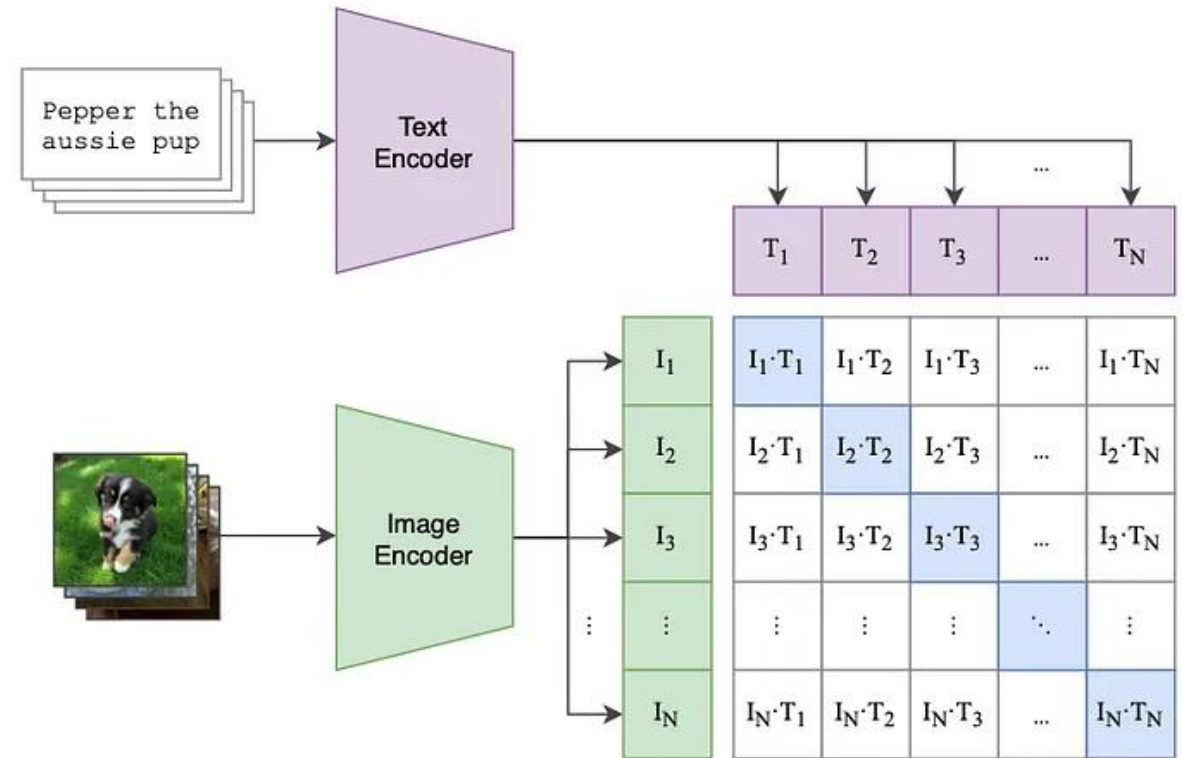
# Demo

Stable Diffusion



# Contrastive Language-Image Pretraining (CLIP)

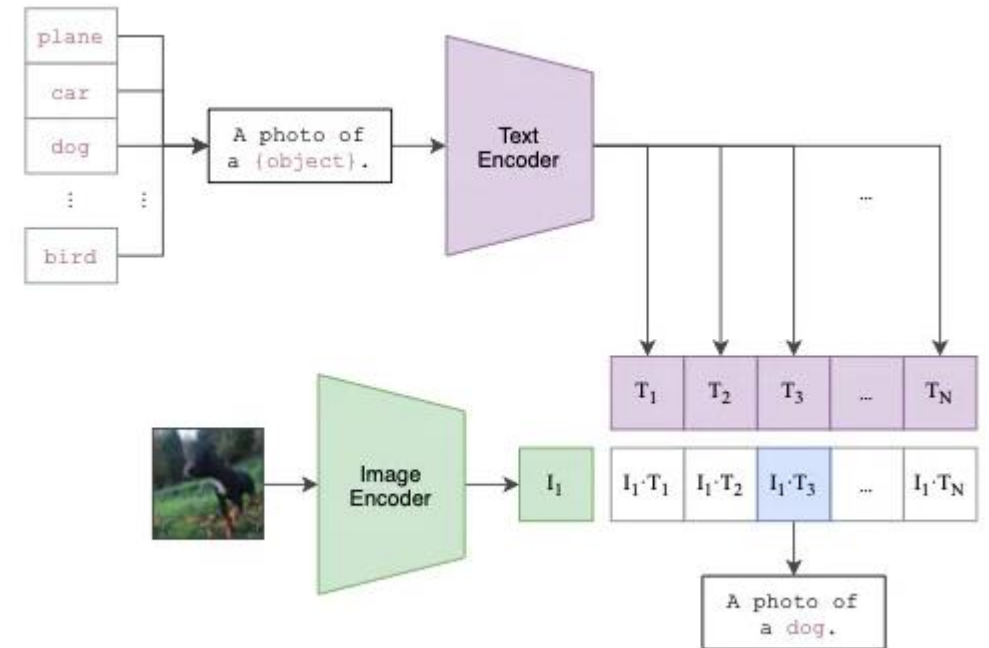
- Introduced in 2021 by OpenAI and trained on more than **400 million** text-image pairs
- Correlates text and images by training a model to maximize embedding similarities
- DALL·E 2 uses CLIP's text encoder to generate embeddings from text prompts





# Zero-Shot Image Classification

- Classify photos without:
  - Training a CNN from scratch or using transfer learning with a pretrained CNN
  - Assembling a labeled dataset
- Present model with an image and several possible descriptions
- Model predicts which description is "correct" by computing similarity of image embedding and each text embedding



# OpenCLIP

- Open-source version of CLIP
  - Hosted at [https://github.com/mlfoundations/open\\_clip](https://github.com/mlfoundations/open_clip)
- OpenAI shared model weights but not the training dataset
- LAION (Large-scale Artificial Intelligence Network) assembled massive text-image datasets, trained several versions of OpenCLIP with them, and published the datasets and model weights
  - Includes LAION-5B dataset with **5.85 billion** text-image pairs
- Many OpenCLIP models trained on these datasets available in Hugging Face's **transformers** package

# Using clip-vit-large-patch14

```
from PIL import Image
from transformers import pipeline

model = pipeline(
    model='openai/clip-vit-large-patch14',
    task='zero-shot-image-classification'
)

image = Image.open('PATH_TO_IMAGE')
model(image, candidate_labels=['owl', 'giraffe', 'camel'])
```

```
[{'score': 0.9978225231170654, 'label': 'giraffe'},
 {'score': 0.002148183062672615, 'label': 'camel'},
 {'score': 2.921208033512812e-05, 'label': 'owl'}]
```

