# Vision Models

Philipp Koehn

25 April 2024

# Image Generation

- How does this work?
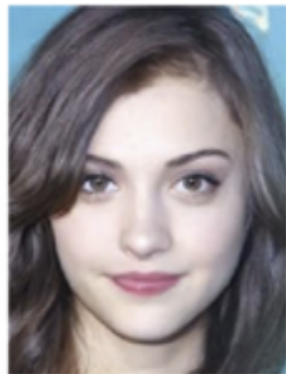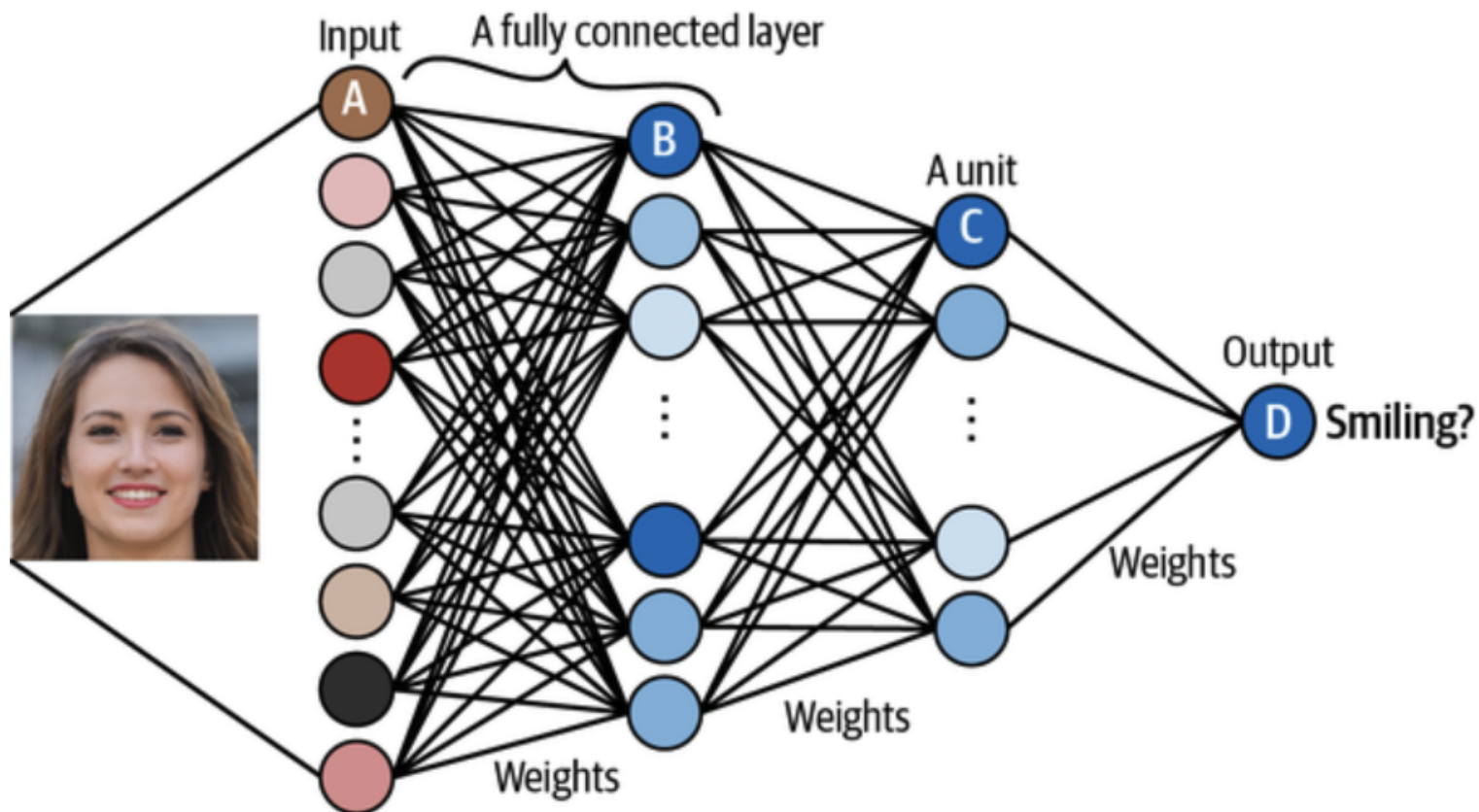
# image classification

# Image Classification

- Input: Image, Output: Class

- 60,000 images

- 32x32 resolution

- 10 classes

# Simple Convolution

- Passing a 3x3 kernel over a 5x5 image

- Using padding to preserve size of representation

# Edge Detection

- Detecting large local color shifts: edges in image

- Two kernels: 2-dimensional vector for each point

# ImageNet: 2012



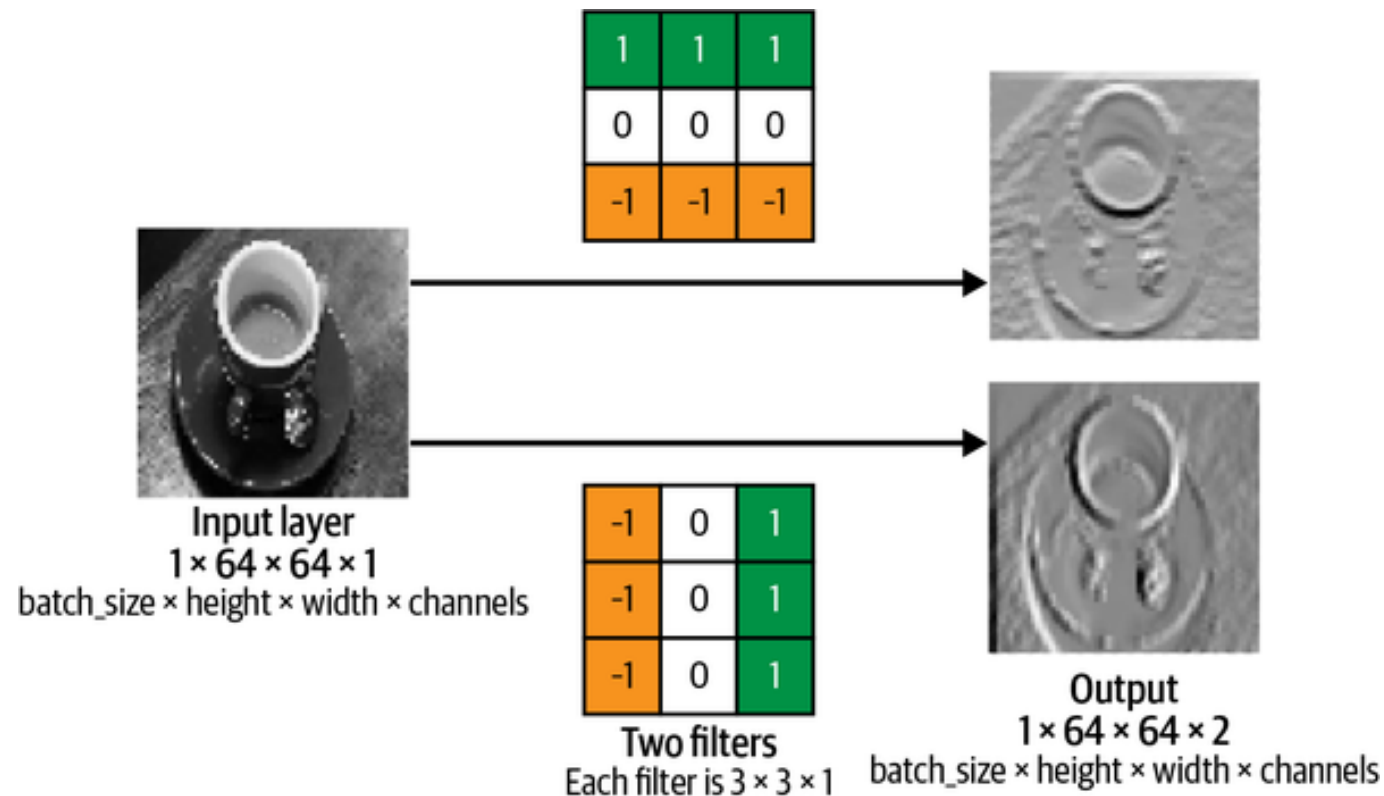- ImageNet: Large scale competition for image classification
  (1.2 million labeled training examples, 1000 categories)

- Breakthrough for deep learning in 2012: large gains with CNN (AlexNet)

# autoencoders

# Encoder-Decoder Model



Example: Image to 2-dimensional vector

# Generating Novel Items of Clothing

- Encoder predicts a normal distribution, specifically the **mean** $\mu$ and **variance** $\sigma^2$

$$f\left(x \mid \mu, \sigma^2\right) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

# Variational Autoencoder

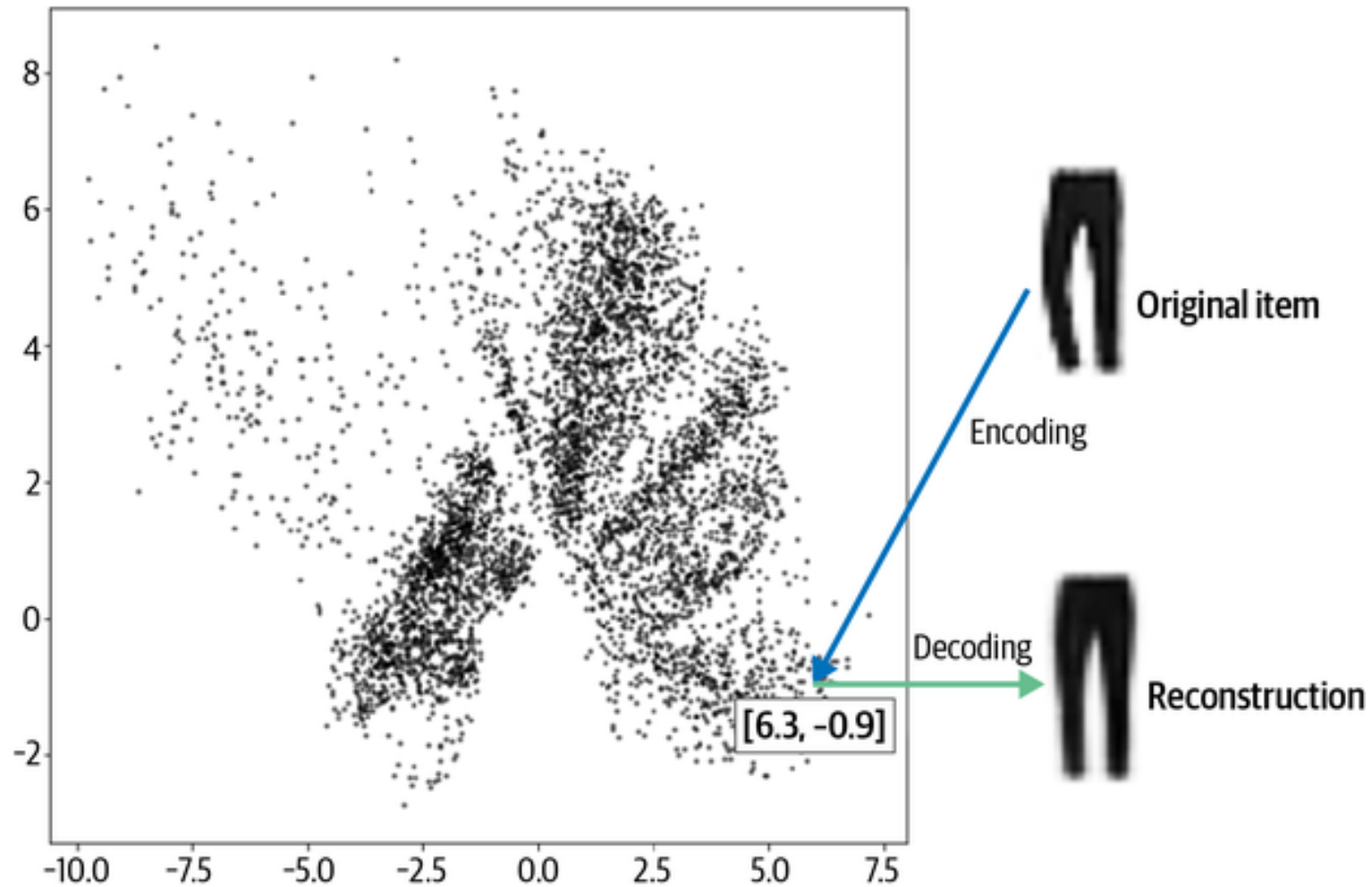- Encoder (as before)

- Predicts mean `z_mean` and variance `z_log_var` of the distribution
  (note: predict 2-dimensional point $\rightarrow$ 2 means and 2 variances)

- We randomly draw a point `z` from the distribution

- Decoder predicts from that randomly drawn point `z` (as before)

- Add preference that predicted mean and variance are normal distribution with mean 0 and variance 1: $N(0, 1)$

- Computed as KL-divergence between predicted values $\mu$ and $\sigma^2$

$$\mathcal{D}_{KL}\left[N(\mu, \sigma) \| N(0, 1)\right] = -\frac{1}{2}\sum\left(1 + \log\left(\sigma^2\right) - \mu^2 - \sigma^2\right)$$

# What is it Good For?

- General principle: adding noise is good (here: random sampling of point z)

- Better use of embedding space (bias towards center)



Left: z_mean, Right: p-values

- CelebA dataset of over 200,000 color images of celebrity faces

- We need a larger model (multiple convolutional layers, larger embedding sizes)



Example real faces

Reconstructions

# Generation of Faces

- Sample points (from a standard multivariate normal distribution)

- Decode the sampled points

- Plot the images

Smiling

- We want to obtain a vector for one feature, e.g. `Smiling`

- Recall: encoder predicts vector representation for each image

- Vector that points in the direction of `Smiling`

  - take average of all vectors for images labelled `Smiling`: $a_{\mathrm{smiling}}$
  - take average of all vectors for images not labelled `Smiling`: $a_{\overline{\mathrm{smiling}}}$
  - subtract the two vectors: $v_{\mathrm{smiling}} = a_{\mathrm{smiling}} - a_{\overline{\mathrm{smiling}}}$

# Generation with Feature Vector

- Randomly sample point $z$ (or use representation for existing image)

- Subtract and feature vector $v_{\text{smiling}}$ with varying factor, generate

- Take representations of two images: $z_A$ and $z_B$

- Combine them with weight $\alpha \in [0, 1]$ into new vector $z = \alpha z_A + (1 - \alpha) z_B$

# generative adversarial training

- The story so far: we can generate novel images

- New task: detect if image is generated or real

# Training

- Train initial generator

- Iterate

  - train discriminator
  - train generator with additional discriminator loss

- Note: when training generator, leave discriminator parameters fixed
  (and vice versa)

# Common Problems

- Discriminator overpowers generator

- Generator overpowers discriminator
  (for instance mode collapse: generates unique image that fools discriminator)

- Training loss of generator is not informative
  (mainly battles discriminator and not fitting the training data)

- Many hyperparameters

- Wasserstein GAN

- Lipschitz Constraint

- Gradient Penalty Loss

- Conditional GAN

# diffusion models

# Diffusion Models

- Key idea

  - create noise image
  - generate image from noise
  - repeat this for, say, 20 steps

- Training data for this process can be creating by adding noise



$$x_0 \qquad x_{t-1} \qquad q(x_t|x_{t-1}) \qquad x_t \qquad x_T$$

- This is done in stages, each time adding Gaussian noise $\epsilon_t$
  (mean 0 and unit variance, but then scaling to effective variance $\beta_t$ below)

$$x_t = \sqrt{1 - \beta_t}\, x_{t-1} + \sqrt{\beta_t}\, \epsilon_{t-1}$$

- Or, written as a probability distribution from one image to the next

$$q\left(x_t \mid x_{t-1}\right) = \mathcal{N}\left(x_t; \sqrt{1 - \beta_t}\, x_{t-1}, \beta_t\, \mathbf{I}\right)$$

- After a large number of steps, this becomes indistinguishable from a noise image

# Diffusion Schemes

- The variance $\beta_t$ is changed throughout the process

  - small changes to initial, original image
  - larger changes towards the end to ensure randomness

- Common options: linear, cosine, offset cosine

# Effect of Diffusion Schemes

- Cosine diffusion scheme makes less changes initially

- Learn a model $p_\theta$ that maps a noisy image $x_t$ back to a less noisy image $x_{t-1}$

- We actually learn a model that maps back $x_t$ to the original image $x_o$
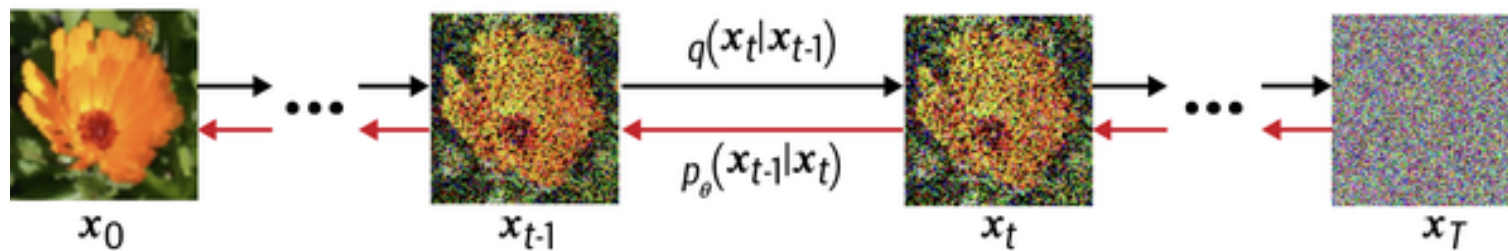
- The accumated noise can be computed
$$q\left(\boldsymbol{x}_t \mid \boldsymbol{x}_0\right) = \mathcal{N}\left(\boldsymbol{x}_t; \sqrt{\bar{\alpha}_t}\,\boldsymbol{x}_0, \left(1 - \bar{\alpha}_t\right)\mathbf{I}\right)$$
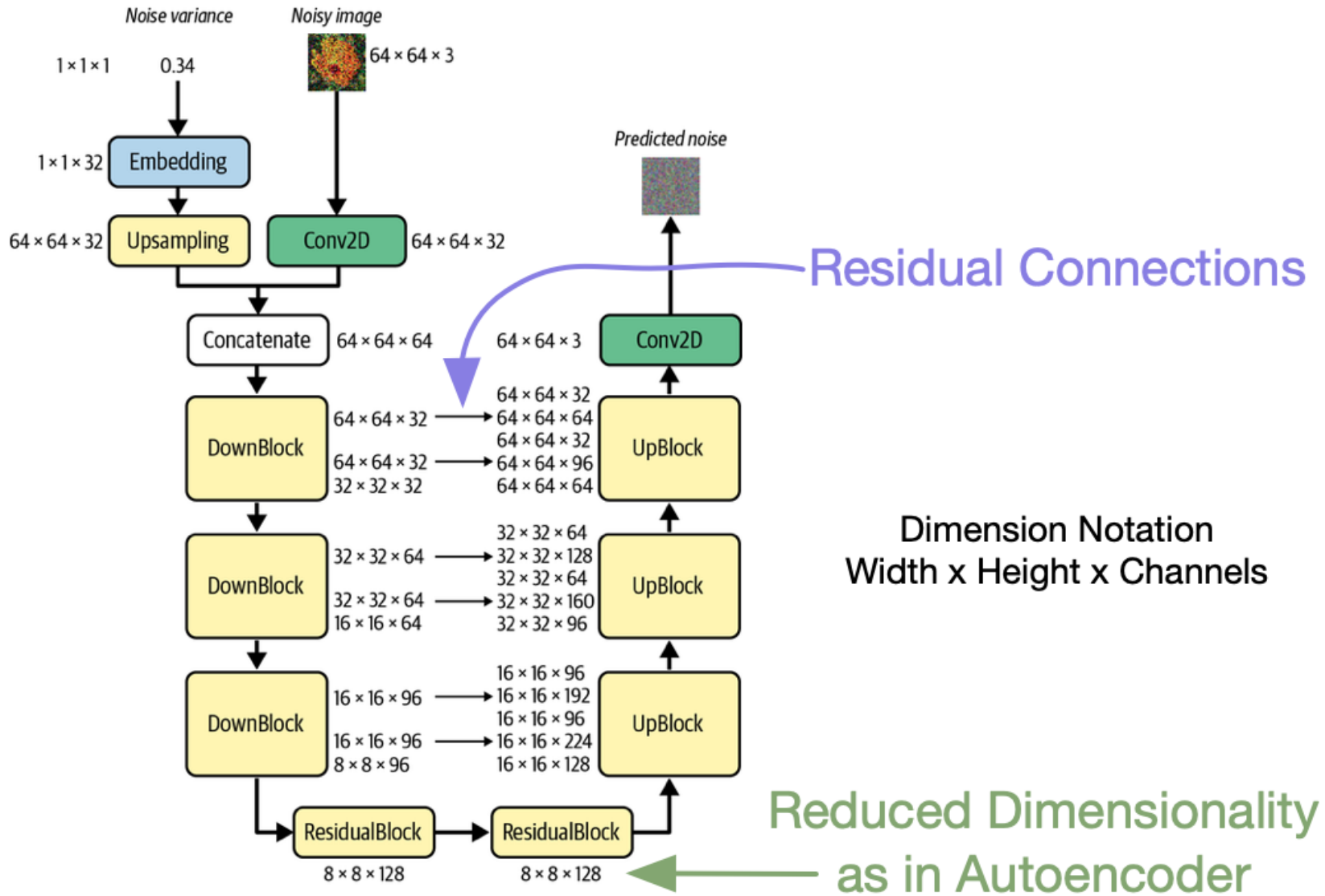
- The model predicts the noise $\epsilon_\theta(x_t)$

**Algorithm 1** Training

1: **repeat**
2: $\quad \mathbf{x}_0 \sim q(\mathbf{x}_0)$
3: $\quad t \sim \text{Uniform}(\{1, \ldots, T\})$
4: $\quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
5: $\quad$ Take gradient descent step on
$$\nabla_\theta \left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(\sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon}, t) \right\|^2$$
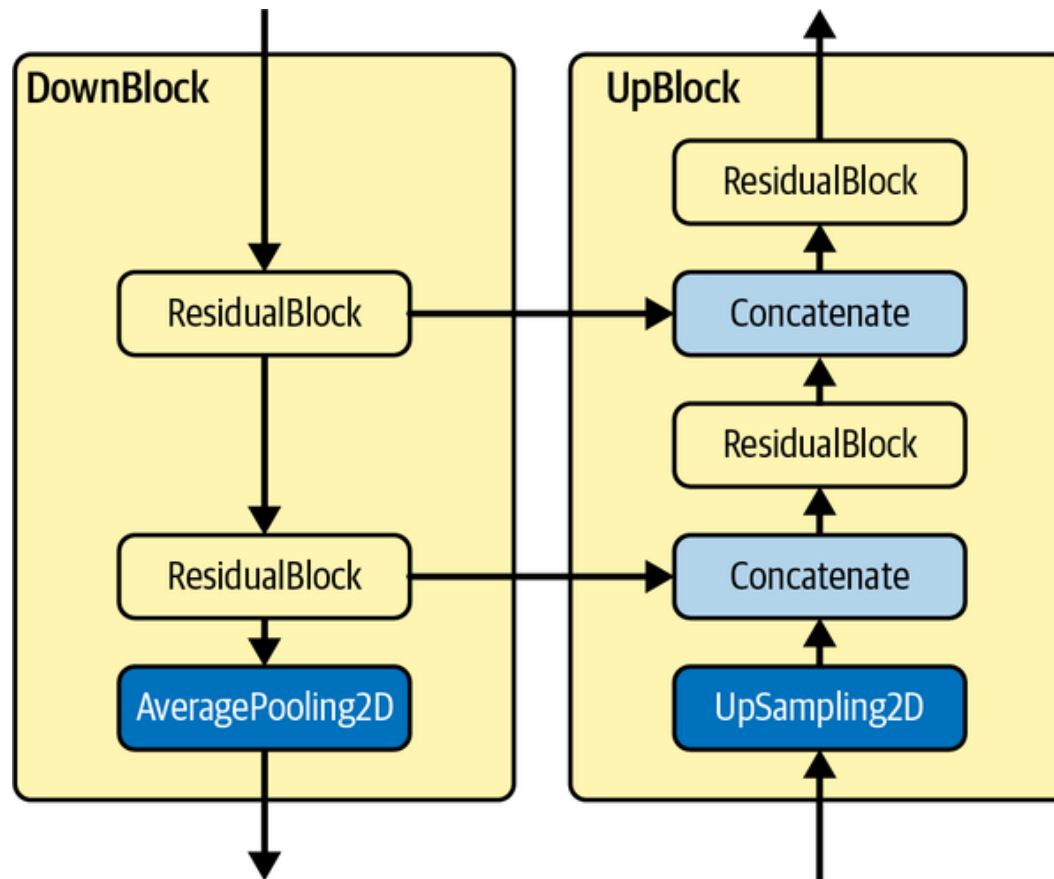6: **until** converged

# Residual Block

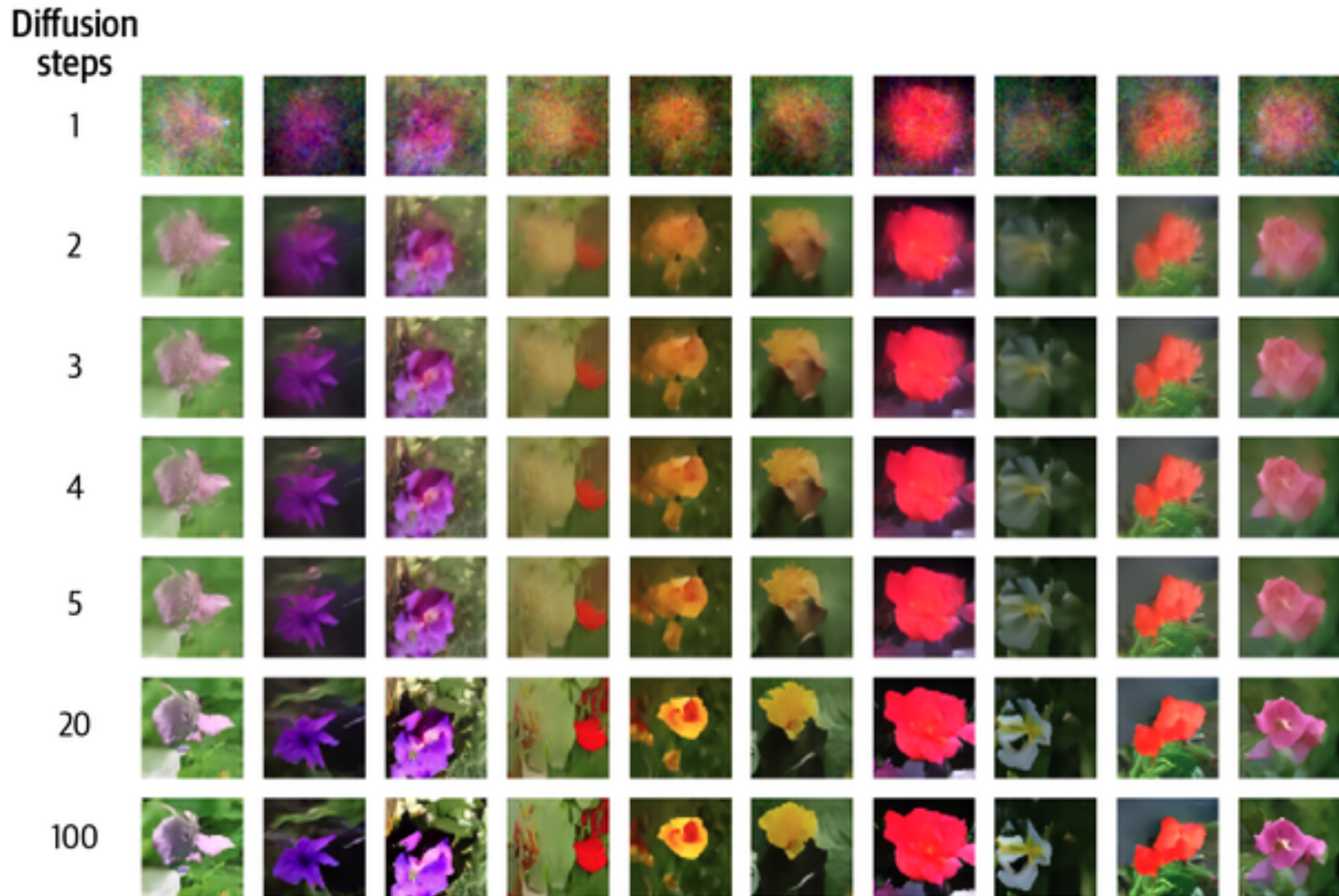(may add convolution with kernel size 1 to residual connection
to generate tensor with the right number of channels)

# DownBlock and UpBlock

# Generation in Multiple Steps



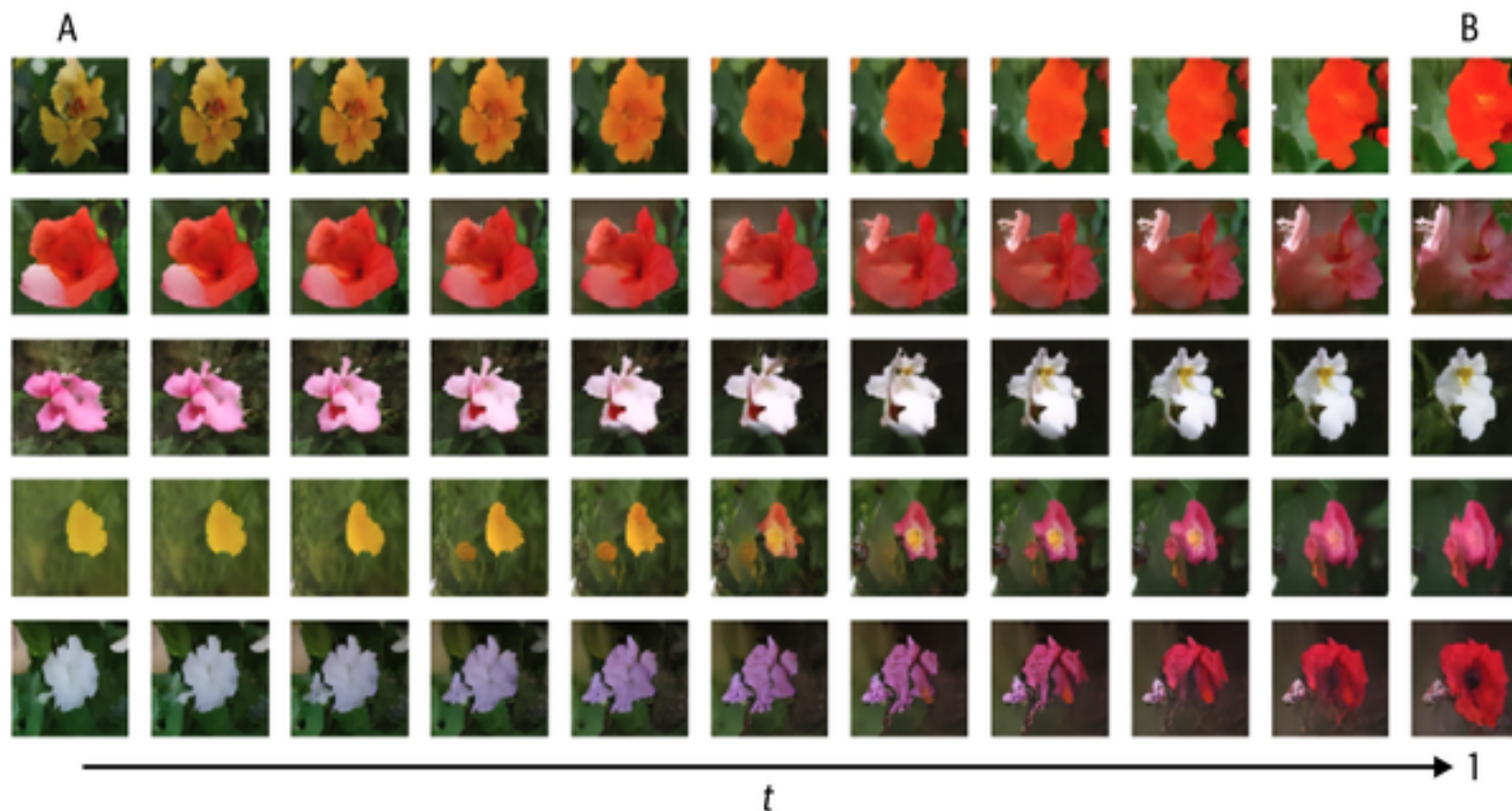**Diffusion steps**

1, 2, 3, 4, 5, 20, 100

# Interpolation

- Generation is deterministic given noise tensor $a$ and $b$

$\Rightarrow$ We can interpolate between different noise images

# Common Failures

# advanced gan

- First train a model for lower resolution images

- Upscale the image in stages

# StyleGAN: Mixing Image Styles

# VQ-GAN: Vector Quantization

- Input: high dimensional vector

- K-Means Clustering



- Centroid vectors form a codebook, each vector is replaced with cluster ID

- Vectors are replaced by their nearest centroid

- Instead of using ConvNNs, Transformer model predicts sequence of patches
- This particular model also uses vector quantization

# text to image generation

A head of broccoli made out of modeling clay, smiling in the sun

"Leonardo da Vinci early sketches of a cyborg"
Input prompt

Text encoder

Text embedding

Prior

Image embedding

Decoder

LENO

Generated image

- Text is encoded as a prior to the image generation process

- Key training step: Contrastive Language-Image Pre-training (CLIP)

- Given pairs of images and text (scraped from the Internet)

- Learn representations of text (using Transformer models)

- Learn representations of images (using ViT-VQ GAN)

- Learn mapping between them

- Representation of text and image as a single vector, mapped to same size

- Training: minimize cosine similarity of real pairs

- Note: this is not a generative model

- Image decoder is a Transformer model

- Predicts patches of the image at each step

- Generation is also conditioned
  on the text representation (the prior)

- Alternatively: diffusion decoder

# stable diffusion

Robach et al. (2022):
High-Resolution Image Synthesis with Latent Diffusion Models

Generated from prompt "a photograph of an astronaut riding a horse"

# Latent Diffusion Model

# Autoencoder

- Convert image into lower-dimensional latent space

- Trained independently as first step

- Refinements
  - KL loss (as in variational autoencoders)
  - patch loss (for better detail rendering)

# Diffusion Model



Operates in latent space — otherwise the same U-Net from before

Attention (Query, Key, Value) to representations of semantic maps, text, images

- External input $y$ is converted into an intermediate representation $\tau_\theta \in \mathbb{R}^{M \times d_r}$

- For use in the attention model, the intermediate representations of the U-Net are also flattened to $\varphi_i(z_t) \in \mathbb{R}^{N \times d_\epsilon^i}$
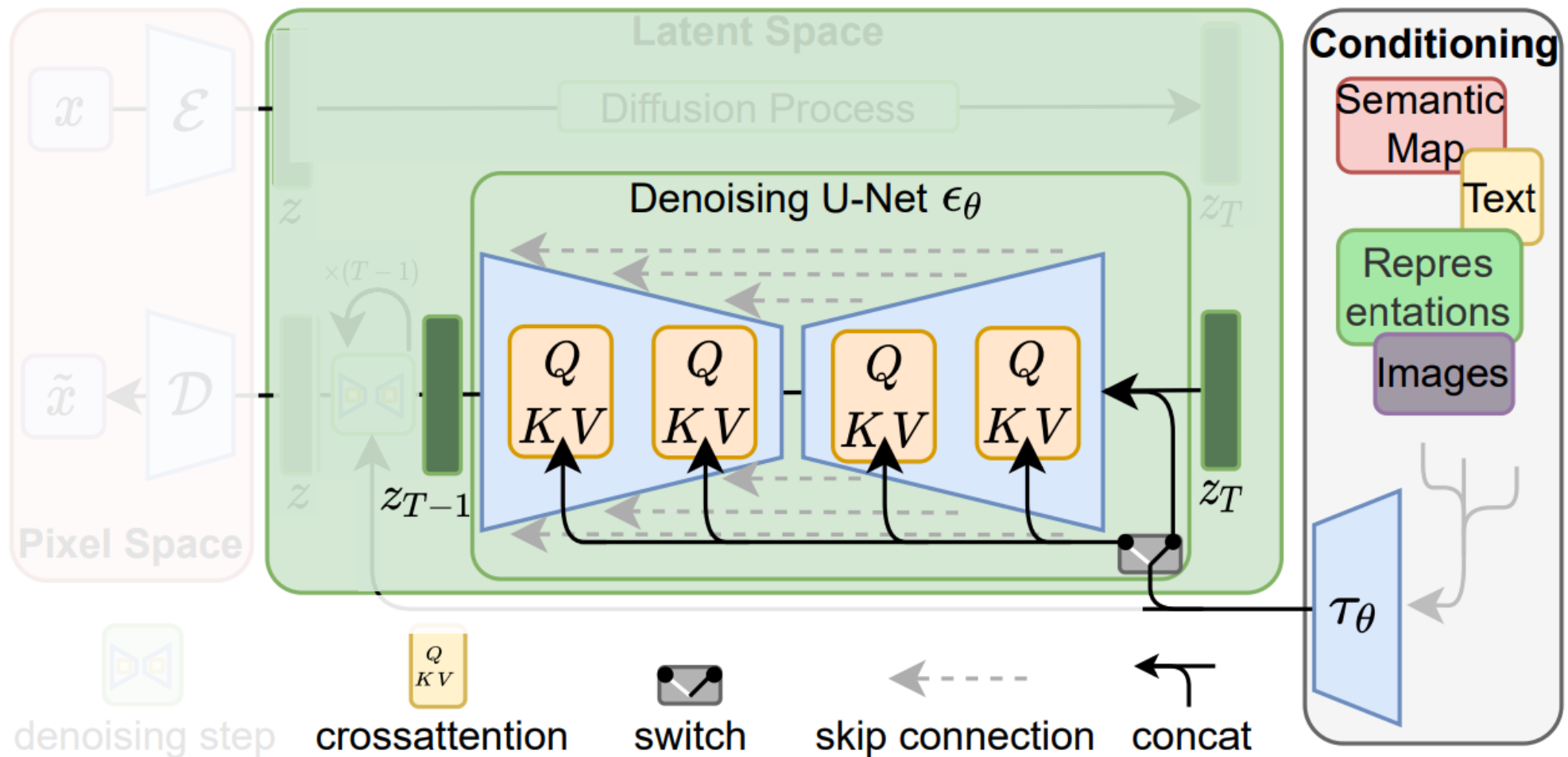
- Attention is
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right) \cdot V$$

  with $\quad Q = W_Q^{(i)} \cdot \varphi_i(z_t) \quad K = W_K^{(i)} \cdot \tau_\theta(y) \quad V = W_V^{(i)} \cdot \tau_\theta(y)$

  that map to vectors of size $d$

- Parameters for $\tau_\theta$ and $\epsilon_\theta$ are jointly optimized using diffusion objective

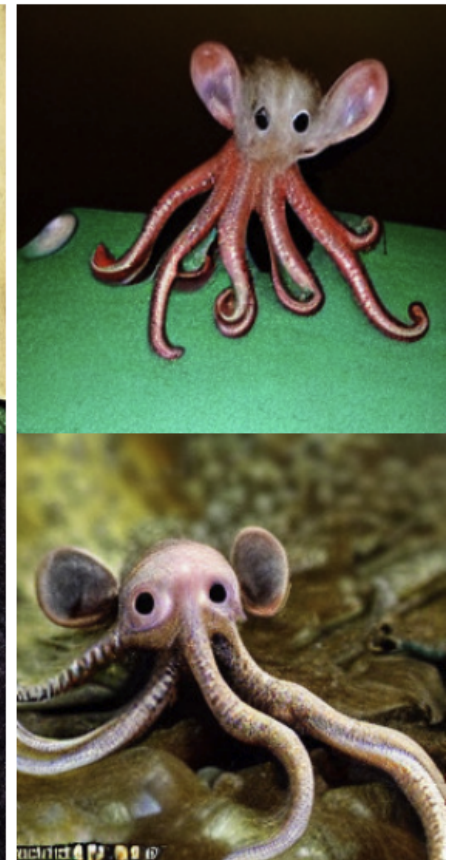'A zombie in the style of Picasso'    'An image of an animal half mouse half octopus'

- Text is represented as a sequence of words

- Transformer model generates $\tau_\theta \in \mathbb{R}^{M \times d_r}$

- Trained on language prompts

  - LAION-400M
  - 400 million text-image pairs
  - extracted from web pages with alt-text in HTML image tag
  - filtered in various ways

- Trained on Open Images dataset of images with labelled object detection

  – 9.2 million images collected from Flickr
  – bounding boxes with object labels
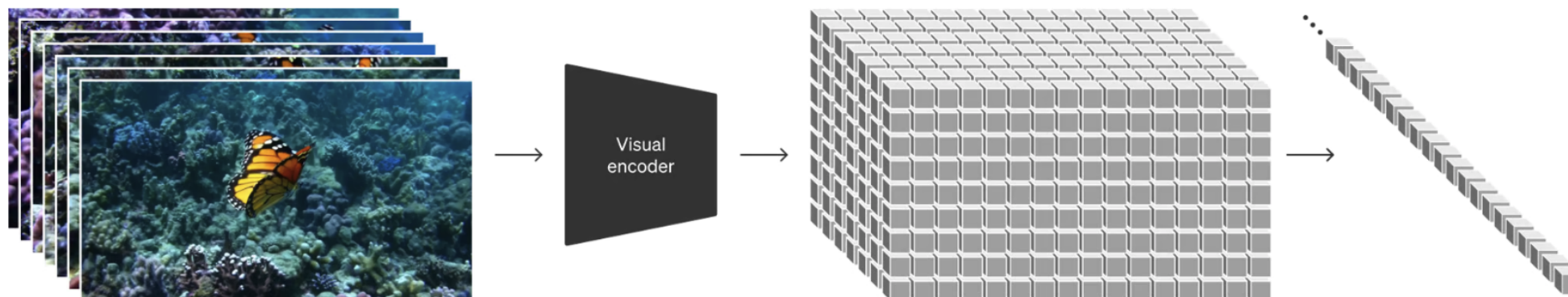  – computer-assisted annotation: automatic labels vetted by humans

- Automatic generation of training data with synthetic masks

- Training aims to reconstruct the original image

# video generation

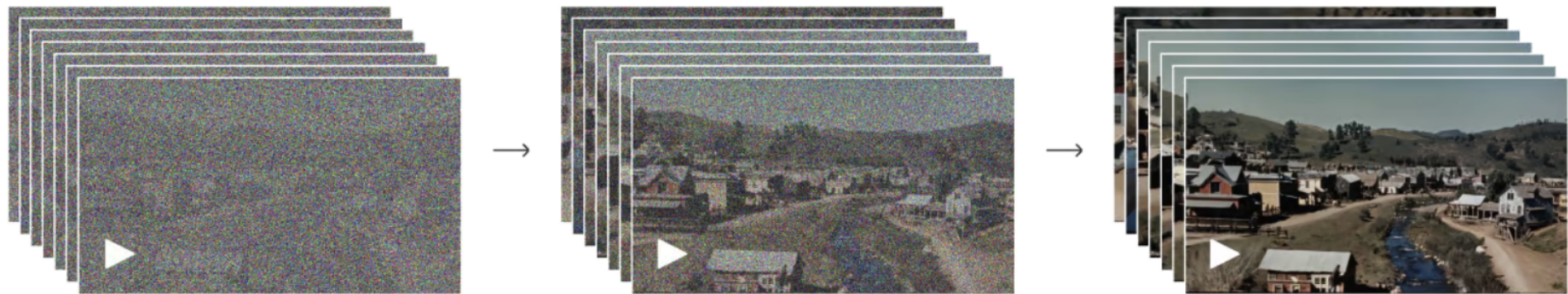OpenAI (2024): Video generation models as world simulators (Sora)

`https://openai.com/research/video-generation-models-as-world-simulators`

- Up to a full minute of high definition video

- Compress image into lower dimensional latent space

- Decompose representation into spacetime patches

- Given noisy patches and conditioning text prompts

- Predict original clean patches using Transformer model

- Training data generated by re-captioning existing videos
  - first, train captioning model
  - use it to produce highly descriptive text captions for video

- Prompting with images (maybe images generated by DALL-E)

- Extending existing video (forward or backward)

- Video-to-video editing (video + text prompt $\rightarrow$ new video)

- Connecting videos

# questions?