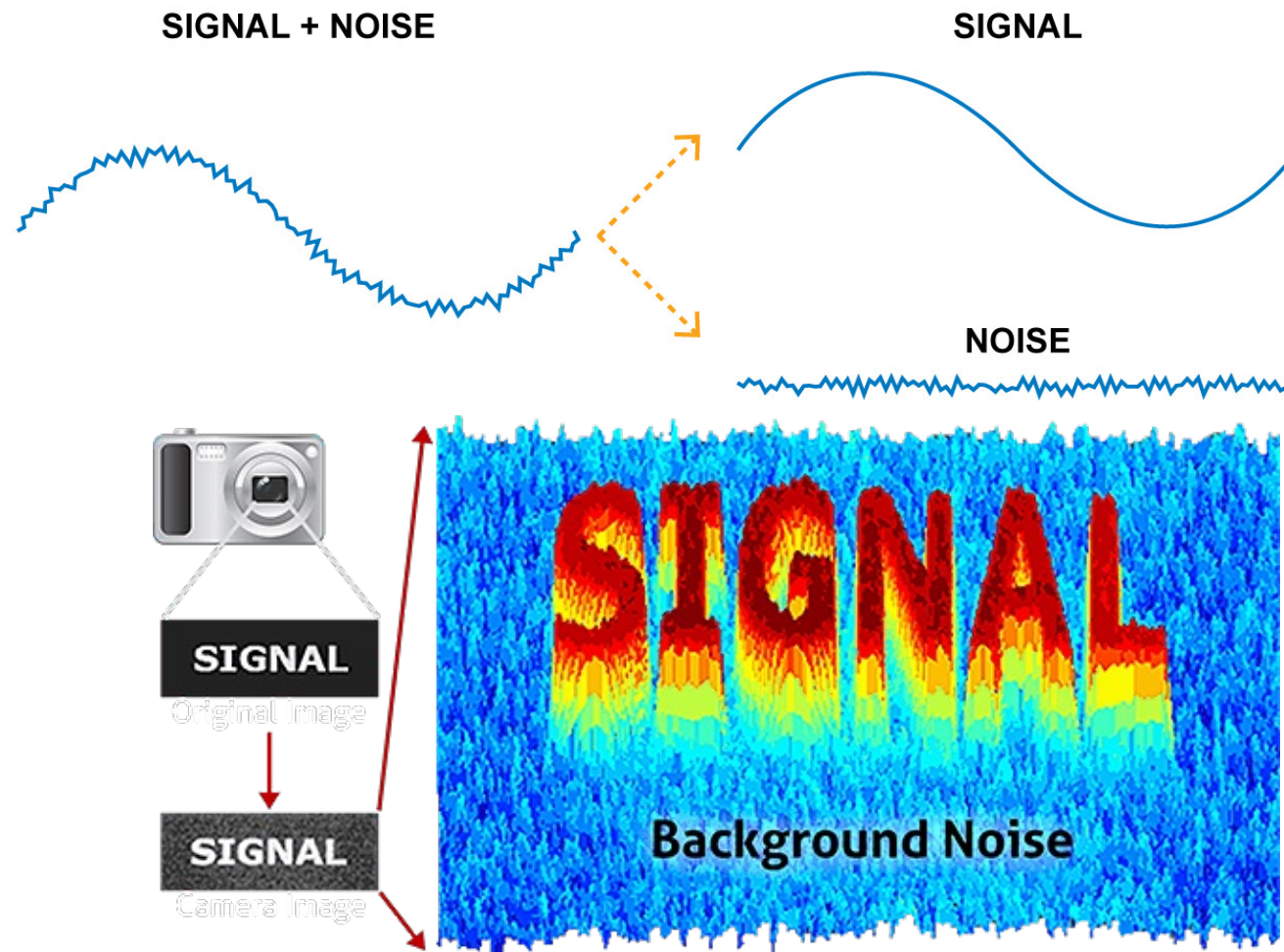# Noise, Pattern, and Image

Noise, Perception, and Learning: Applications in AI Art
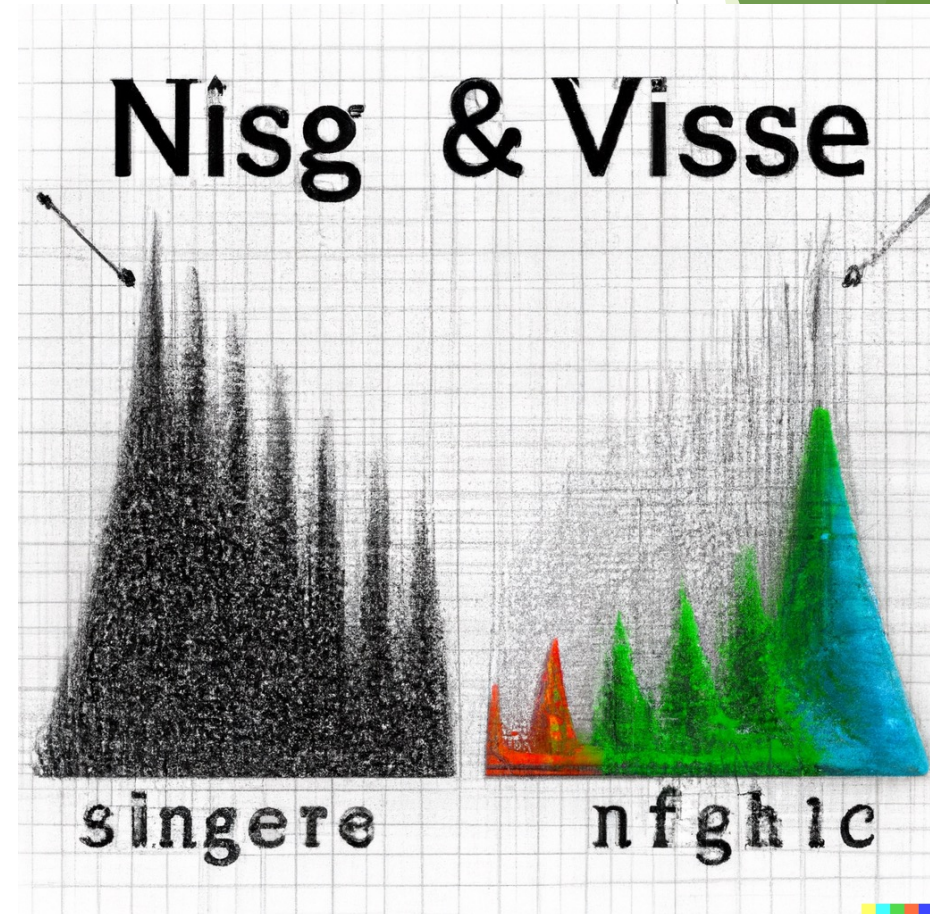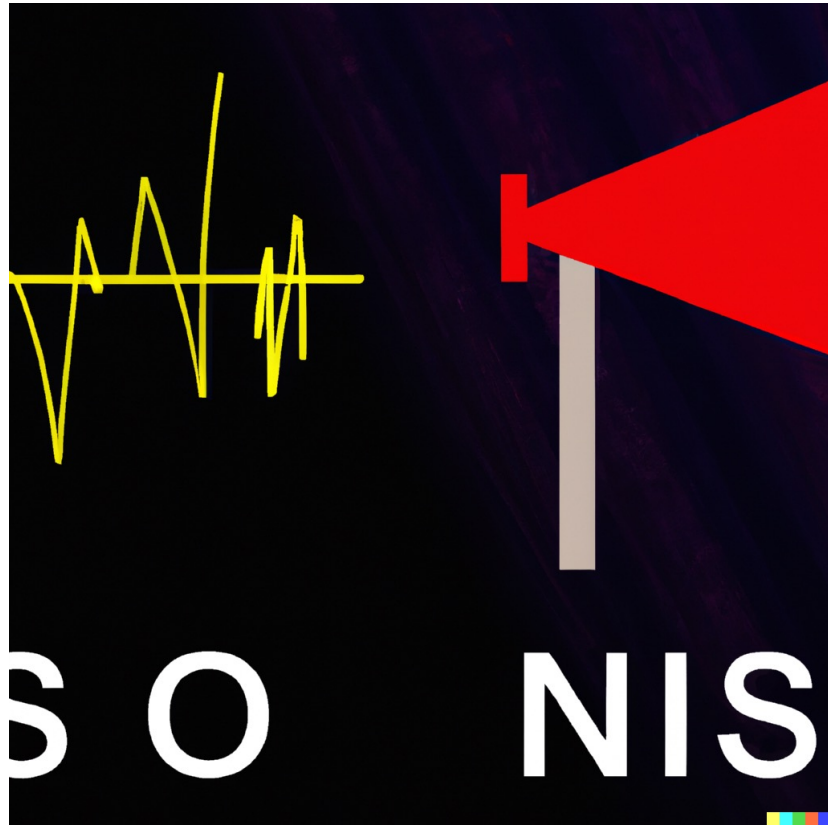
IAP 2023

Sarah  Muschinske
01/25/2023

# What is Noise?



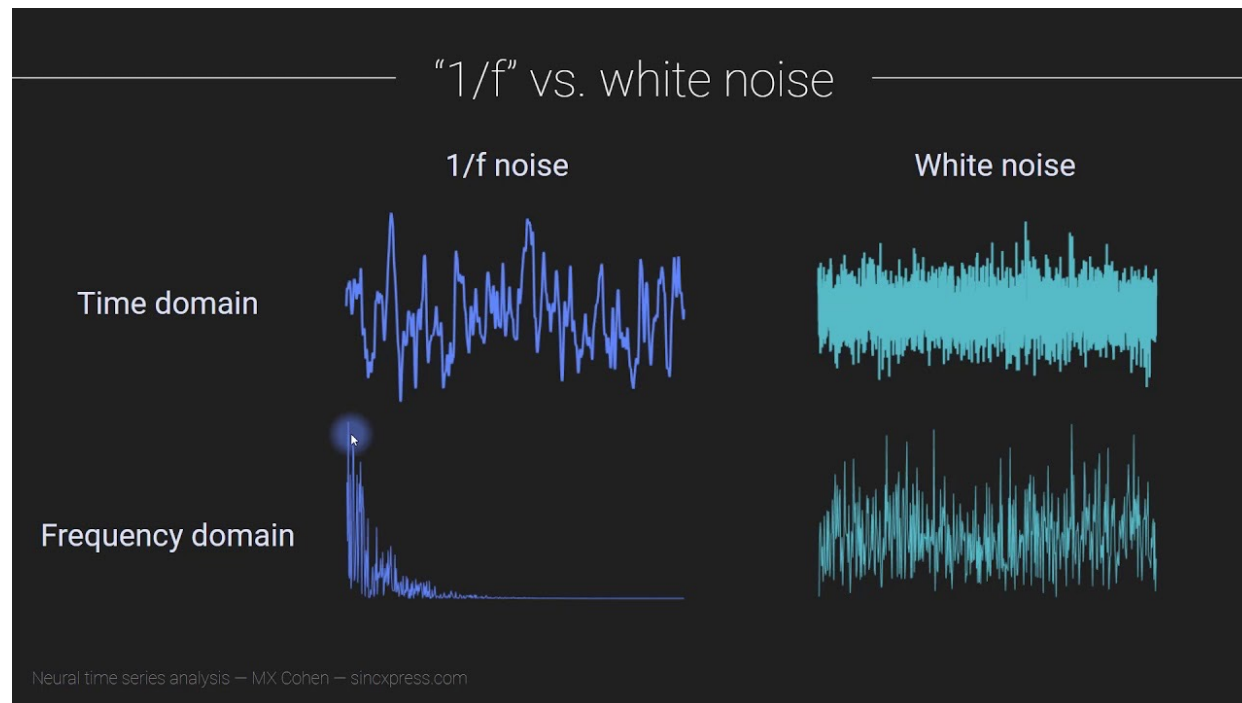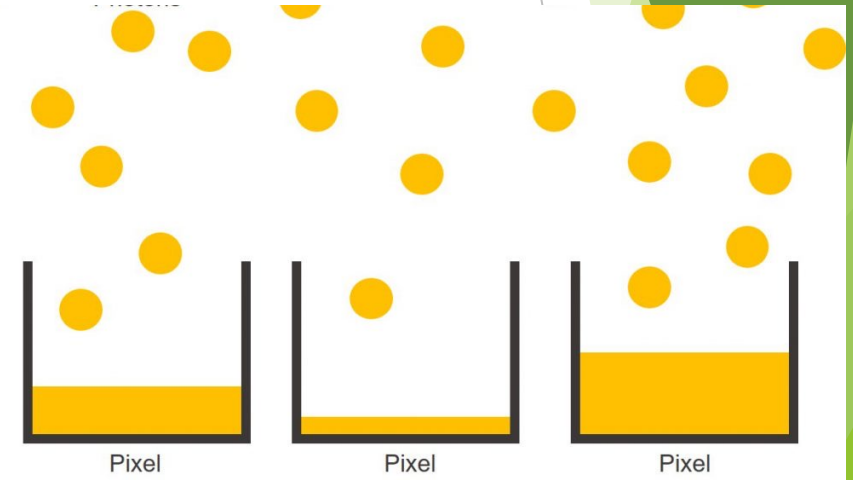SIGNAL + NOISE

SIGNAL

NOISE

SIGNAL
Original Image

SIGNAL
Camera Image

Background Noise

# What does DALL-E 2 think signal-to-noise is?

# Quantum noise



"1/f" vs. white noise

1/f noise — White noise

Time domain

Frequency domain

Neural time series analysis — MX Cohen — sincxpress.com

Photon shot noise



Pixel    Pixel    Pixel
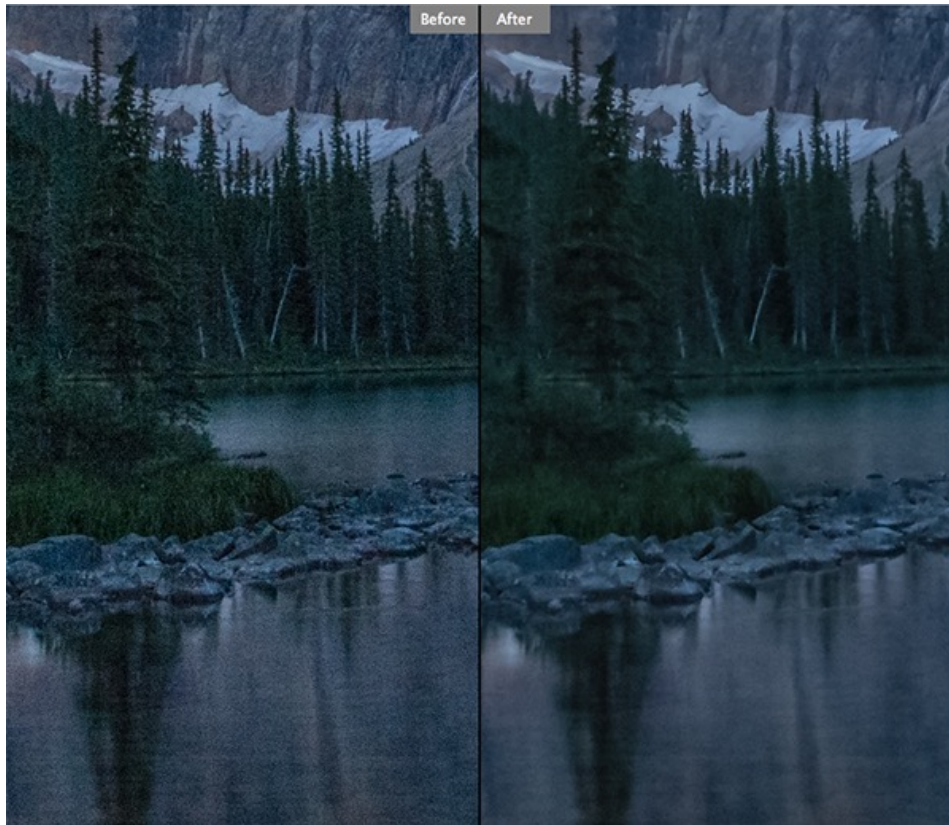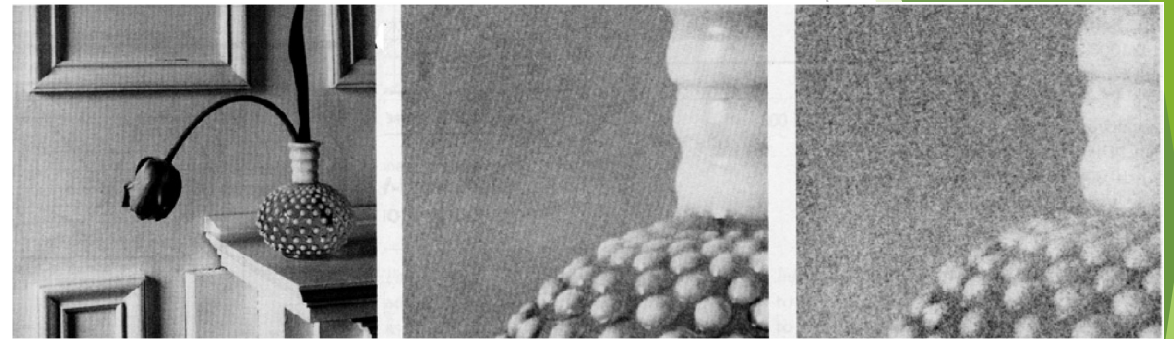
# Noise in Image
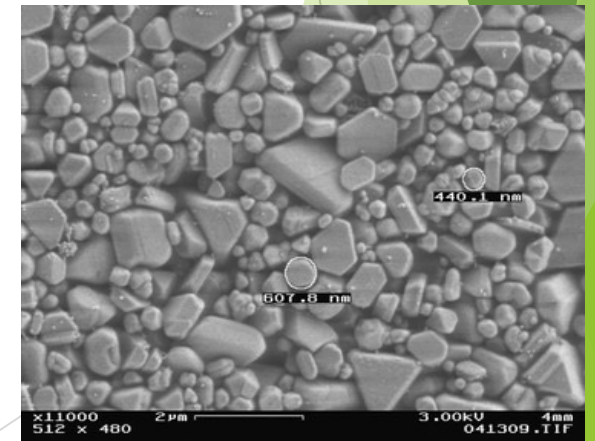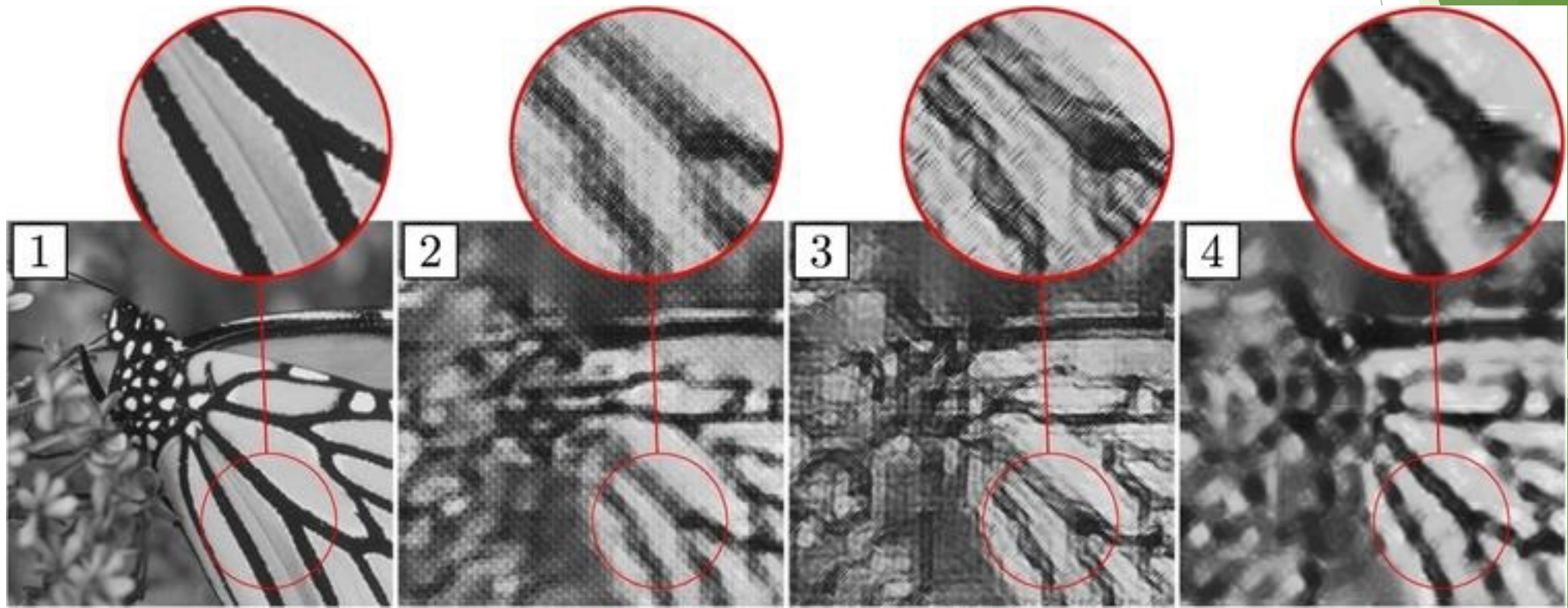
Digital

Film





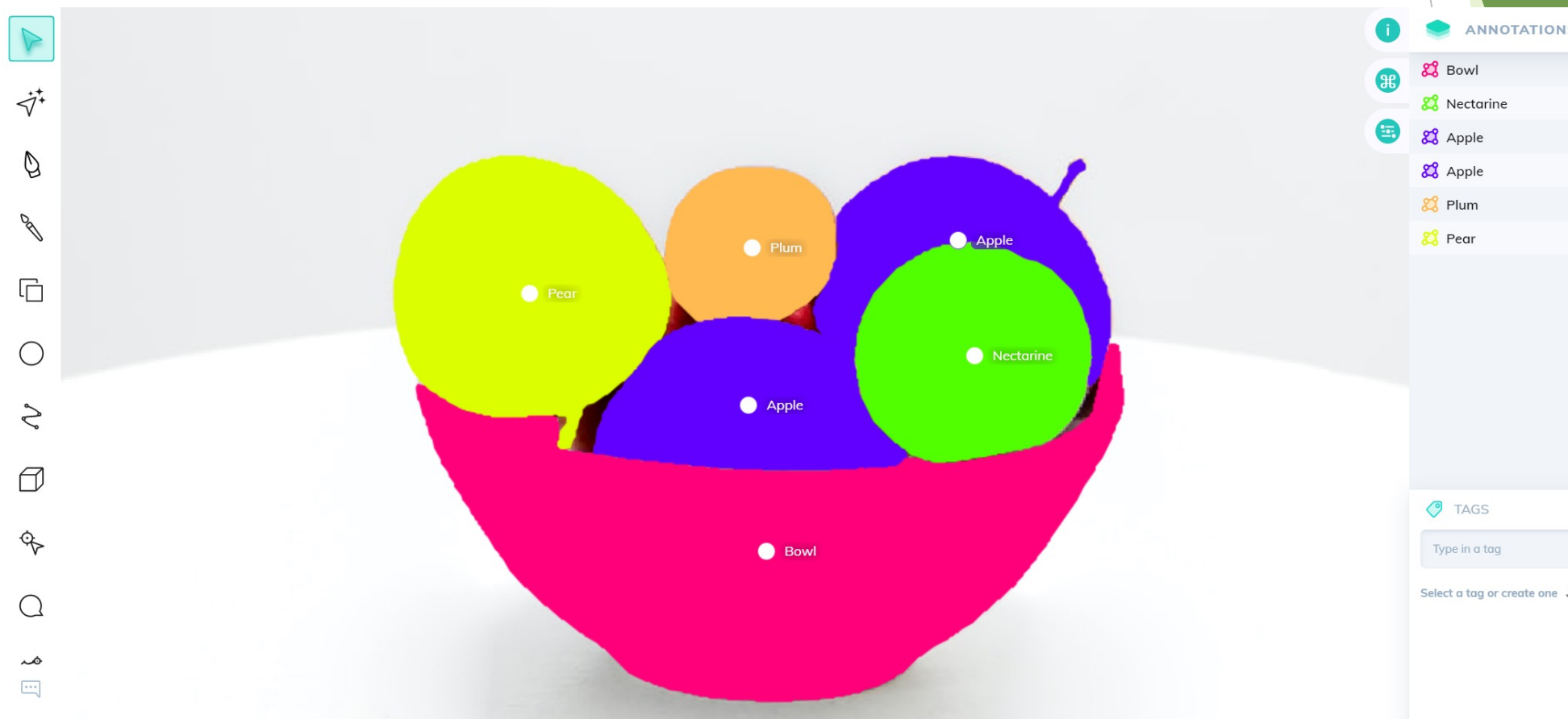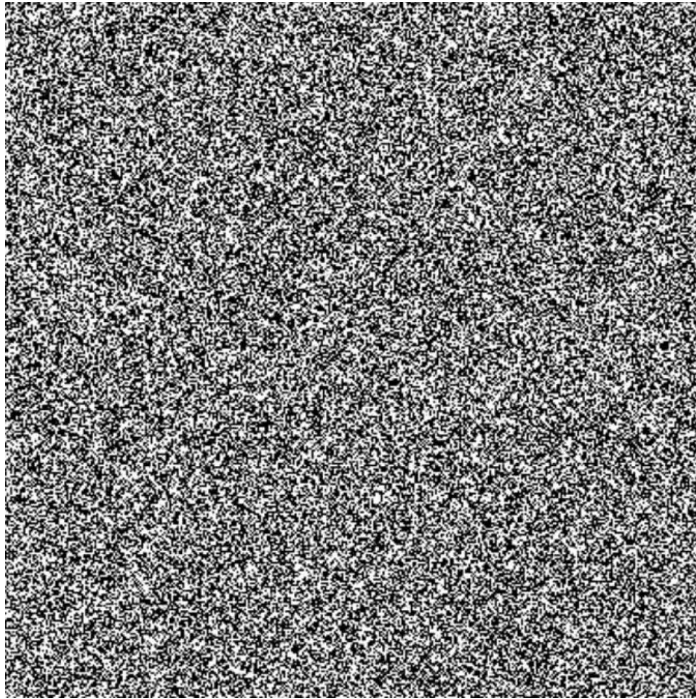Kodak H1 film grain a.) zoomed out b.) fine grain c.) coarse grain


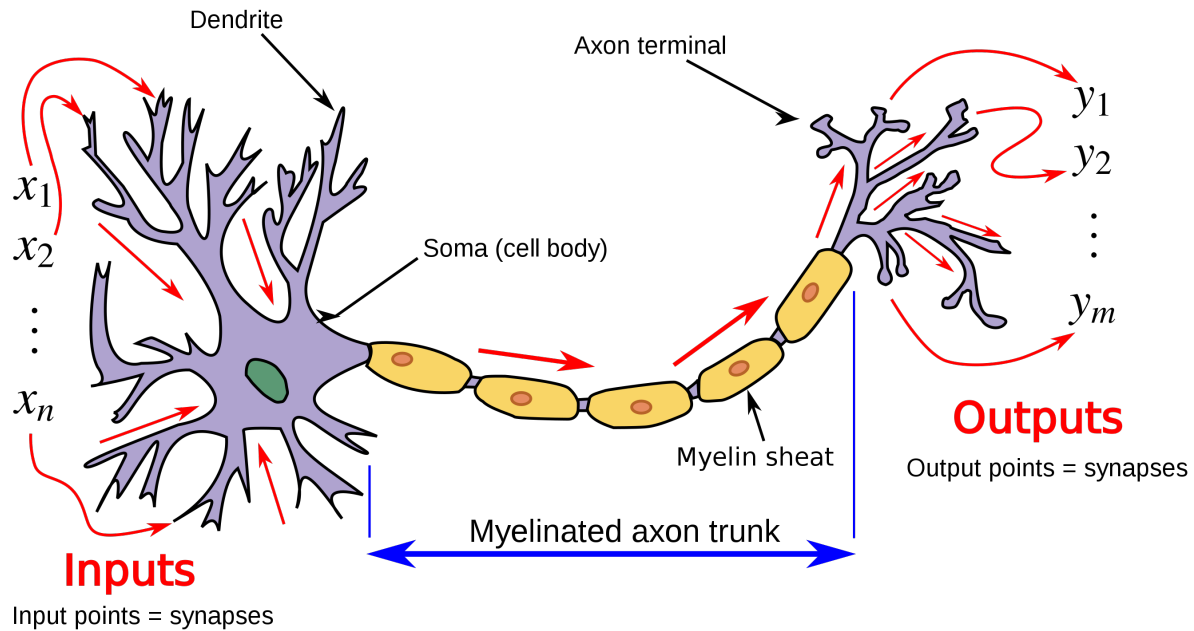
SEM of film grain

# Perceptual Compression

# Semantic Segmentation

# Text-to-Image Generation

# Neural Networks

**Natural**



Dendrite

Axon terminal

$x_1$
$x_2$
$\vdots$
$x_n$

Soma (cell body)

Myelin sheat

**Outputs**

$y_1$
$y_2$
$\vdots$
$y_m$

Output points = synapses

**Myelinated axon trunk**

**Inputs**

Input points = synapses

**Artificial**



Inputs

x0=1

b

x1

w1

x2

w2

wn

...

xn

z  f

y

Output

Linear function

Activation function

# ImageNet

# Wordnet



(a)Semantic hierarchy

entity

animal    vehicle

dog    bird    car    boat

(b)Accuracy of prediction

1

1    0

0    1    0    0

↑ ground truth

(c)Reward of prediction

0

1    0

0    2    0    0

↑ ground truth

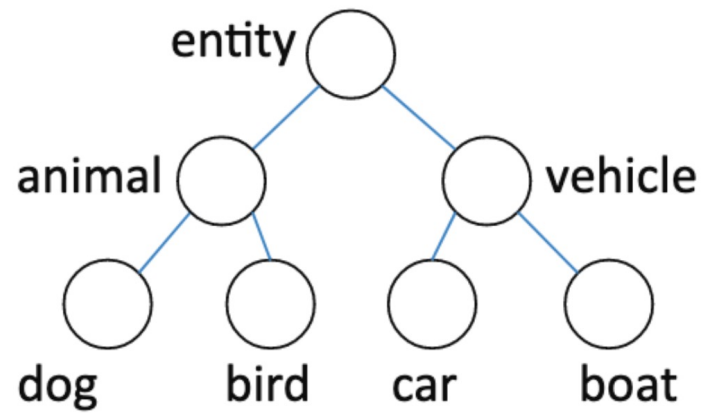# Classifiers

# Style Transfer



$$\hat{y} = \arg \min_{y} \lambda_c \ell_{feat}^{\phi,j}(y, y_c) + \lambda_s \ell_{style}^{\phi,J}(y, y_s) + \lambda_{TV} \ell_{TV}(y)$$

# Super-resolution



|  | **Ground Truth** | **Bicubic** | **Ours** ($\ell_{pixel}$) | **SRCNN** [11] | **Ours** ($\ell_{feat}$) |
|---|---|---|---|---|---|
| This image | | 31.78 / 0.8577 | 31.47 / 0.8573 | 32.99 / 0.8784 | 29.24 / 0.7841 |
| Set5 mean | | 28.43 / 0.8114 | 28.40 / 0.8205 | 30.48 / 0.8628 | 27.09 / 0.7680 |

# GANs

thispersondoesnotexist.com

# This x does not exist

# GAN Play

▶ https://mitmedialab.github.io/GAN-play/



https://phillipi.github.io/pix2pix/ -- the relevant paper 2017.
https://affinelayer.com/pixsrv/
Pix-to-pix

# GAUGAN

▶ http://gaugan.org/gaugan2/

# Stable Diffusion

# Diffusion

- Predicts the score function $\nabla_x \log p(x)$ for an unconditional model

- Adding conditioning:

$$\nabla_x \log p(x|y) = \nabla_x \log p(y|x) + \nabla_x \log p(x)$$

Where y is your conditioning i.e your text input

$$\nabla_x \log p_\gamma(x|y) = \gamma \nabla_x \log p(y|x) + \nabla_x \log p(x)$$

Where $\gamma$ is the guidance scale

$$\nabla_x \log p_\gamma(x|y) = (1 - \gamma)\nabla_x \log p(x) + \gamma \nabla_x \log p(x|y)$$

# VAE (Variational Autoencoder)

# U-Net

# Denoising



Figure 2: The directed graphical model considered in this work.

# Schedulers

- Sets how much noise the decoder tries to remove with each step

# CLiP



(1) Contrastive pre-training

Pepper the aussie pup → Text Encoder

Image Encoder

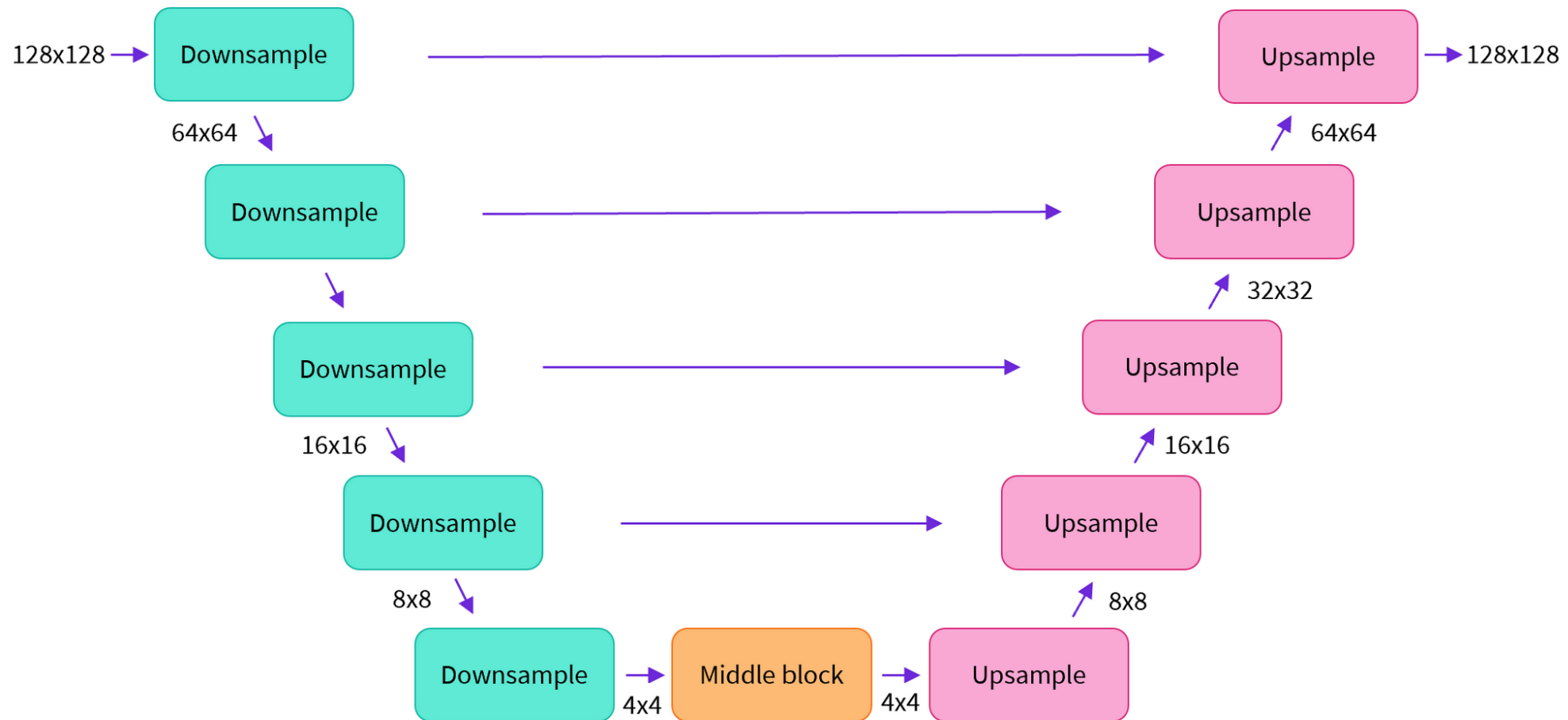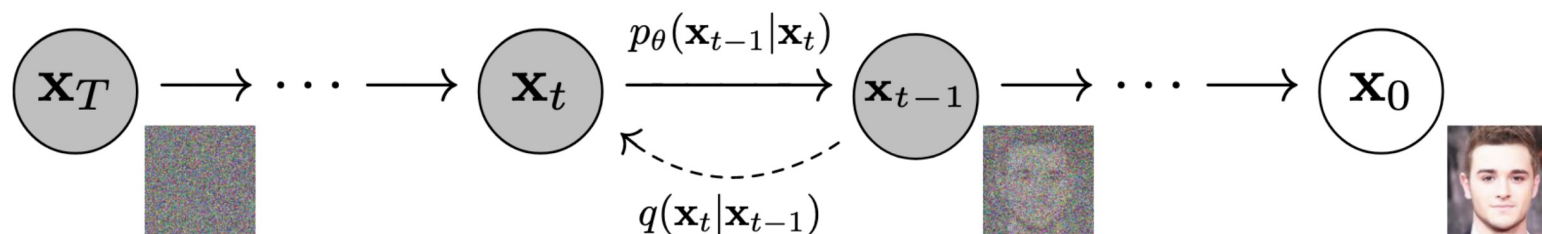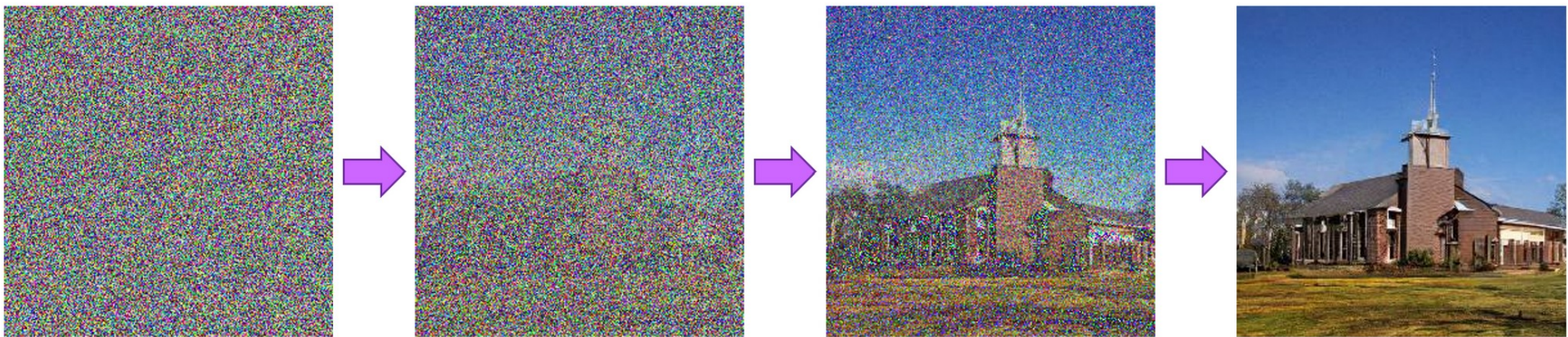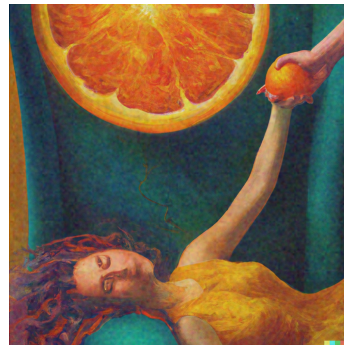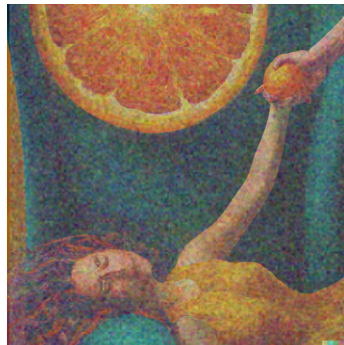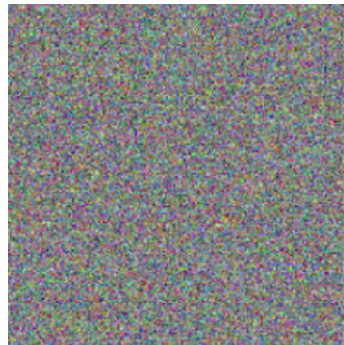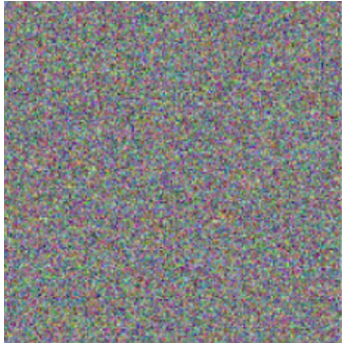| | $T_1$ | $T_2$ | $T_3$ | ... | $T_N$ |
|---|---|---|---|---|---|
| $I_1$ | $I_1 \cdot T_1$ | $I_1 \cdot T_2$ | $I_1 \cdot T_3$ | ... | $I_1 \cdot T_N$ |
| $I_2$ | $I_2 \cdot T_1$ | $I_2 \cdot T_2$ | $I_2 \cdot T_3$ | ... | $I_2 \cdot T_N$ |
| $I_3$ | $I_3 \cdot T_1$ | $I_3 \cdot T_2$ | $I_3 \cdot T_3$ | ... | $I_3 \cdot T_N$ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\ddots$ | $\vdots$ |
| $I_N$ | $I_N \cdot T_1$ | $I_N \cdot T_2$ | $I_N \cdot T_3$ | ... | $I_N \cdot T_N$ |

(2) Create dataset classifier from label text

plane
car
dog
⋮
bird

A photo of a {object}. → Text Encoder

| $T_1$ | $T_2$ | $T_3$ | ... | $T_N$ |
|---|---|---|---|---|

(3) Use for zero-shot prediction

Image Encoder → $I_1$

| $I_1 \cdot T_1$ | $I_1 \cdot T_2$ | $I_1 \cdot T_3$ | ... | $I_1 \cdot T_N$ |
|---|---|---|---|---|

A photo of a dog.

# Ontological Model

# Tokenizer

"Friends, Romans and Countrymen"

friends

romans

countrymen

# Tokenizer

"**Friends, Romans and Countrymen**"

friends

romans

countrymen

# Text Embedding

family friends romans Italy

# Hugging Face



# The AI community building the future.

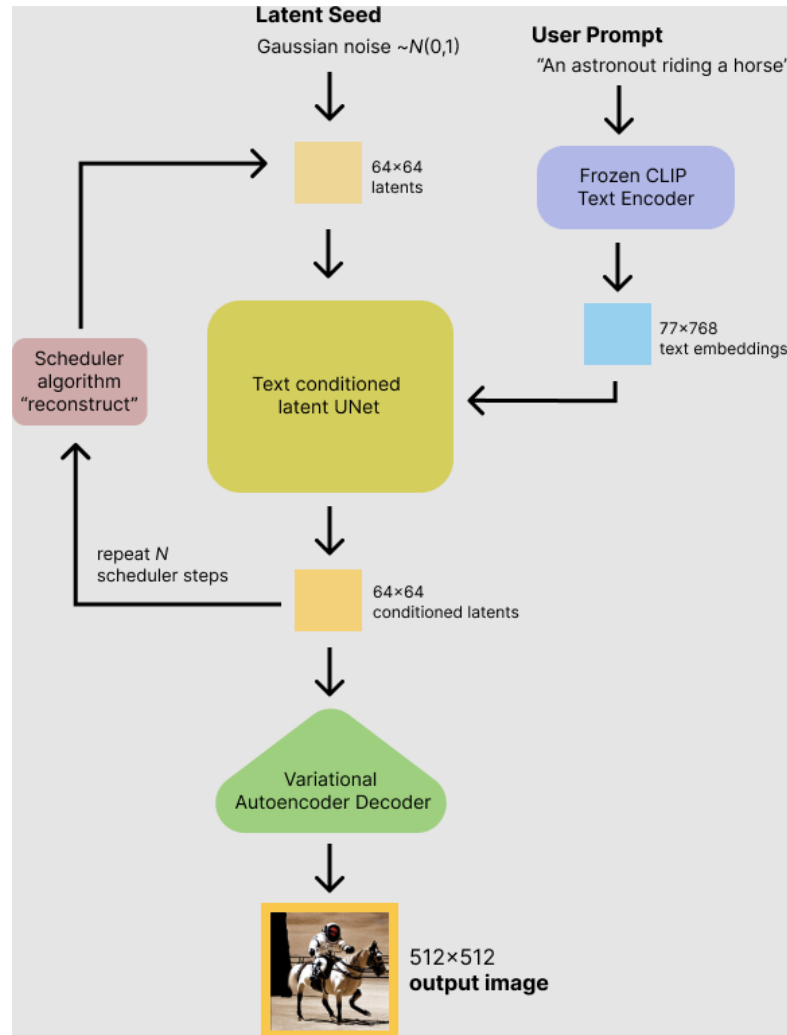Build, train and deploy state of the art models powered by the reference open source in natural language processing.

# Stable Diffusion API with HuggingFace

# Let's try it out

# Imagen

- ▶ Discovered language model trained only on text data are good text encoders for text-to-image

- ▶ Increasing the size of a text-only language model improves output quality more efficiently than increasing the size of an image diffusion model

- ▶ Dynamic thresholding

# Latent Space (mathematically)

▶ items resembling each other are positioned closer to one another in the latent space

▶ Embedding                                              $f: X \rightarrow Y$

   ▶ Def: An instance of a mathematical structure that is contained within another instance such as the rational numbers within integers

   ▶ Must be injective (i.e. 1:1)

   $$f(x_1) = f(x_2) \; implies \; x_1 = x_2$$

# Latent Space (intuitively)