

ROB 498/599: Deep Learning for Robot Perception (DeepRob)

Lecture 21: NeRF; Gaussian Splatting;
Generative Models; Diffusion Model

04/02/2025



<https://deeprob.org/w25/>

Today

- Feedback and Recap (5min)
- Lightning Talk Review (5min)
- Nerf (20min)
- Gaussian Splatting (15min)
- Generative Models (10min)
- Diffusion Models (15min)
- Summary and Takeaways (5min)

Some feedback on the “Lightning Talk”

Fire!



*Comments and feedback will be released via **Canvas**

- 5-min: long or short?
- Connection to lecture content - Thanks! (LORA, pose, RL, ResNet, surface norm, cnn/ViT/nerf/splatting, etc.....)
- Dataset
- Be mindful of training time/timeline in general
- Collaboration - divide and conquer
- Encourage explore more recent SOTA methods
- *Bonus for live robot demo!*
- *Encourage further extension of the work*

April 22nd, final project showcase @FRB atrium

April 28th, final project (report, code, video/[website](#)) DUE

Neural Radiance Field (NeRF) For View Synthesis

<https://arxiv.org/pdf/2003.08934>

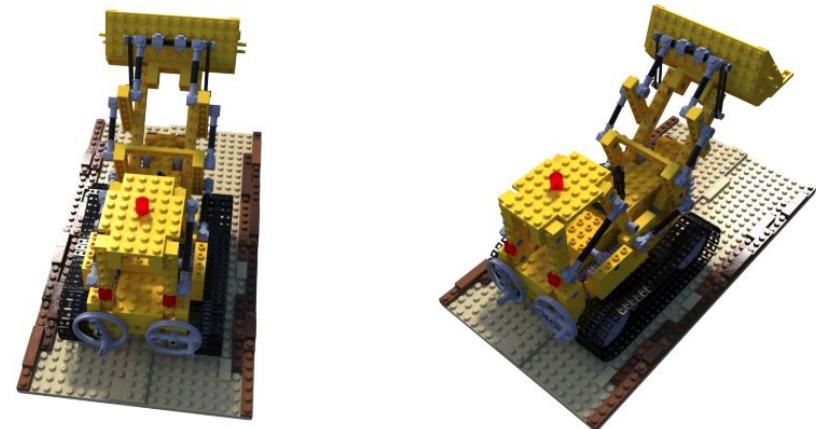
View Synthesis

<https://www.matthewtancik.com/nerf>

Input: Many images of the same scene
(with known camera parameters)

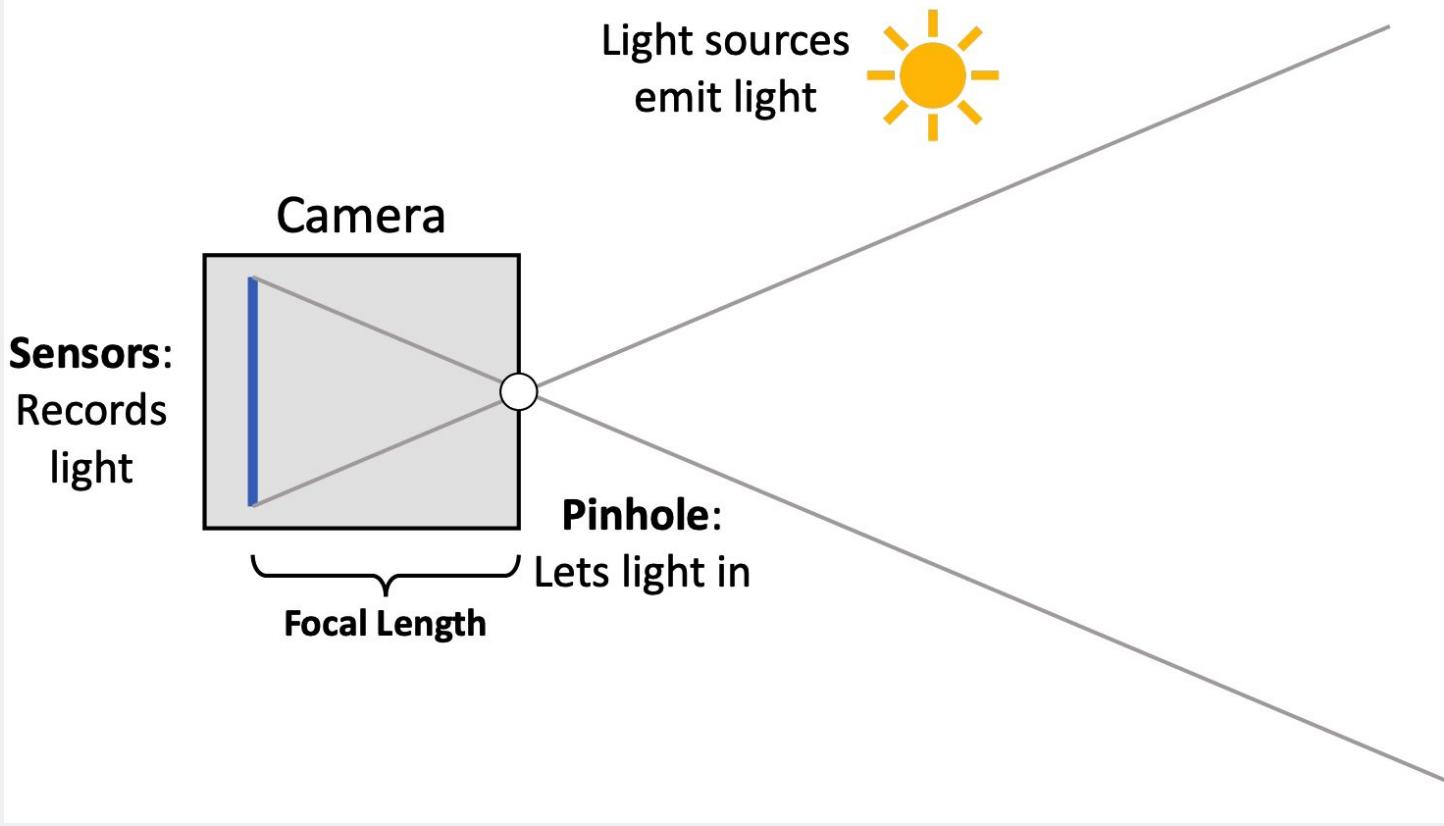


Output: Images showing the scene from novel viewpoints



Scene rendering

Stepping Back: Pinhole Camera Model

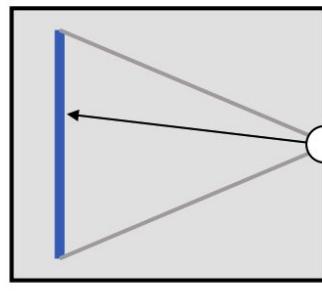


Stepping Back: Pinhole Camera Model

Reflected light
captured by
the camera

Camera

Sensors:
Records
light



Pinhole:
Lets light in

Light sources
emit light



Scenario 1

Objects in the
world reflect light

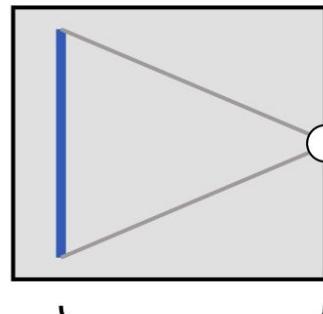


Stepping Back: Pinhole Camera Model

Reflected light
captured by
the camera

Camera

Sensors:
Records
light



Pinhole:
Lets light in

Light sources
emit light



Opaque objects
block light



Scenario 2

Objects in the
world reflect light

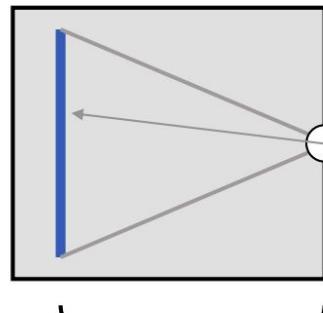


Stepping Back: Pinhole Camera Model

Reflected light
captured by
the camera

Camera

Sensors:
Records
light

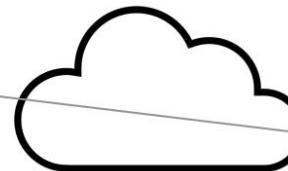


Pinhole:
Lets light in

Light sources
emit light



Transparent objects
attenuate light



Scenario 3

Objects in the
world reflect light

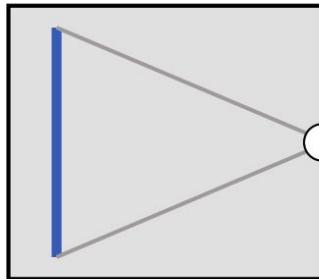


Volume Rendering

Abstract away light sources, objects.

For each point in space, need to know:

- (1) How much light does it emit?
- (2) How opaque is it? $\sigma \in [0,1]$

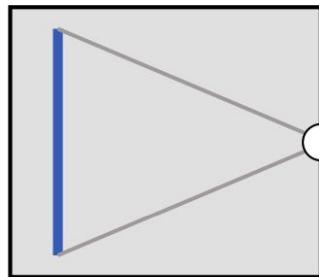


Stepping Back: Pinhole Camera Model

Abstract away light sources, objects.

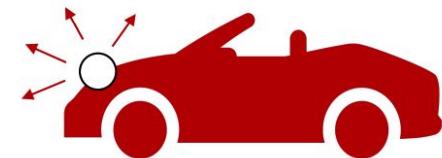
For each point in space, need to know:

- (1) How much light does it emit?
- (2) How opaque is it? $\sigma \in [0,1]$



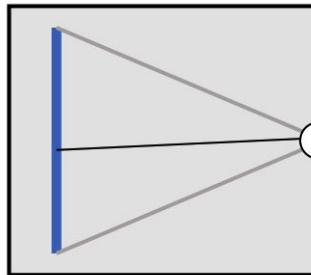
Point in empty space:
(1) Emits no light (black)
(2) Completely transparent $\sigma = 0$

Point on car:
(1) Emits red light in hemisphere
(2) Complete opaque
 $\sigma = 1$



Q: When $\sigma \in (0,1)$?

Volume Rendering



Ray origin

Parameterize each ray as origin plus

direction: $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$

Volume Density is $\sigma(\mathbf{p}) \in [0,1]$

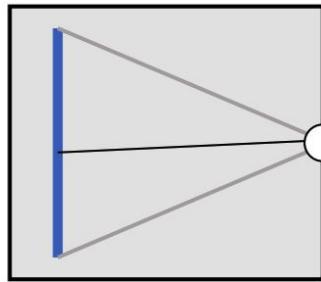
Color that a point \mathbf{p} emits in direction \mathbf{d}

is $c(\mathbf{p}, \mathbf{d}) \in [0,1]^3$

Q: why volume density is designed to be **independent** of the viewing direction, but color depends on viewing direction?

Volume Rendering

Kajiya, J.T., Herzen, B.P.V.: Ray tracing volume densities. Computer Graphics (SIGGRAPH) (1984)



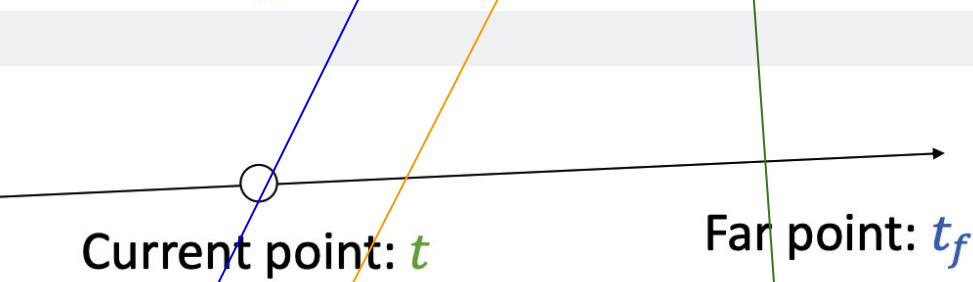
Parameterize each ray as origin plus direction: $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$

Volume Density is $\sigma(\mathbf{p}) \in [0,1]$

Color that a point \mathbf{p} emits in direction \mathbf{d} is $c(\mathbf{p}, \mathbf{d}) \in [0,1]^3$

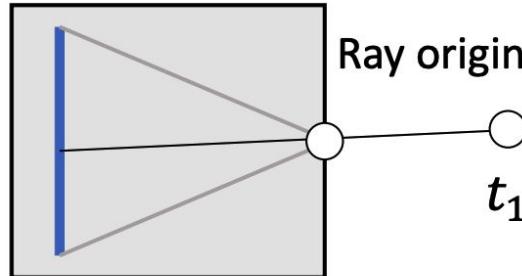
Color observed by the camera given by **volume rendering equation**:

$$C(\mathbf{r}) = \int_{t_n}^{t_f} T(t) \sigma(\mathbf{r}(t)) c(\mathbf{r}(t), \mathbf{d}) dt$$



Transmittance

Compute transmittance by
**accumulating volume
density** up to current point



Parameterize each ray as origin plus
direction: $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$

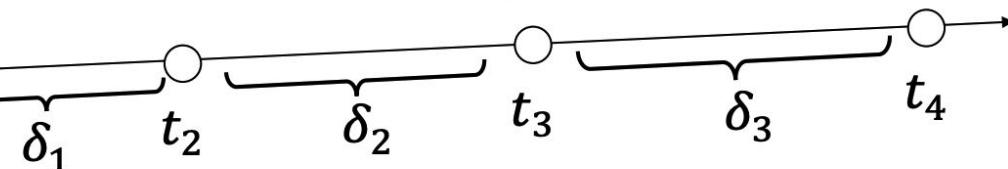
Volume Density is $\sigma(\mathbf{p}) \in [0,1]$

Color that a point \mathbf{p} emits in direction \mathbf{d}
is $\mathbf{c}(\mathbf{p}, \mathbf{d}) \in [0,1]^3$

Color observed by the camera given
by **volume rendering equation**:

$$C(\mathbf{r}) = \int_{t_n}^{t_f} T(t) \sigma(\mathbf{r}(t)) \mathbf{c}(\mathbf{r}(t), \mathbf{d}) dt$$

$$T(t) = \exp \left(- \int_{t_n}^t \sigma(\mathbf{r}(s)) ds \right)$$



Approximate integrals with a set of samples:

$$C(\mathbf{r}) \approx \sum_{i=1}^N T_i (1 - \exp(-\sigma_i \delta_i)) \mathbf{c}_i$$

$$T_i = \exp \left(- \sum_{j=1}^{i-1} \sigma_j \delta_j \right)$$

Mildenhall et al, "Representing
Scenes as Neural Radiance Fields
for View Synthesis", ECCV 2020

NeRF

spatial location (x, y, z)

Train a neural network to input position p and direction d , output $\sigma(p)$ and $c(p, d)$

viewing direction (θ, ϕ)

Training loss: Estimated pixel colors $C(\mathbf{r})$ should match actual pixel colors from images

Volume density

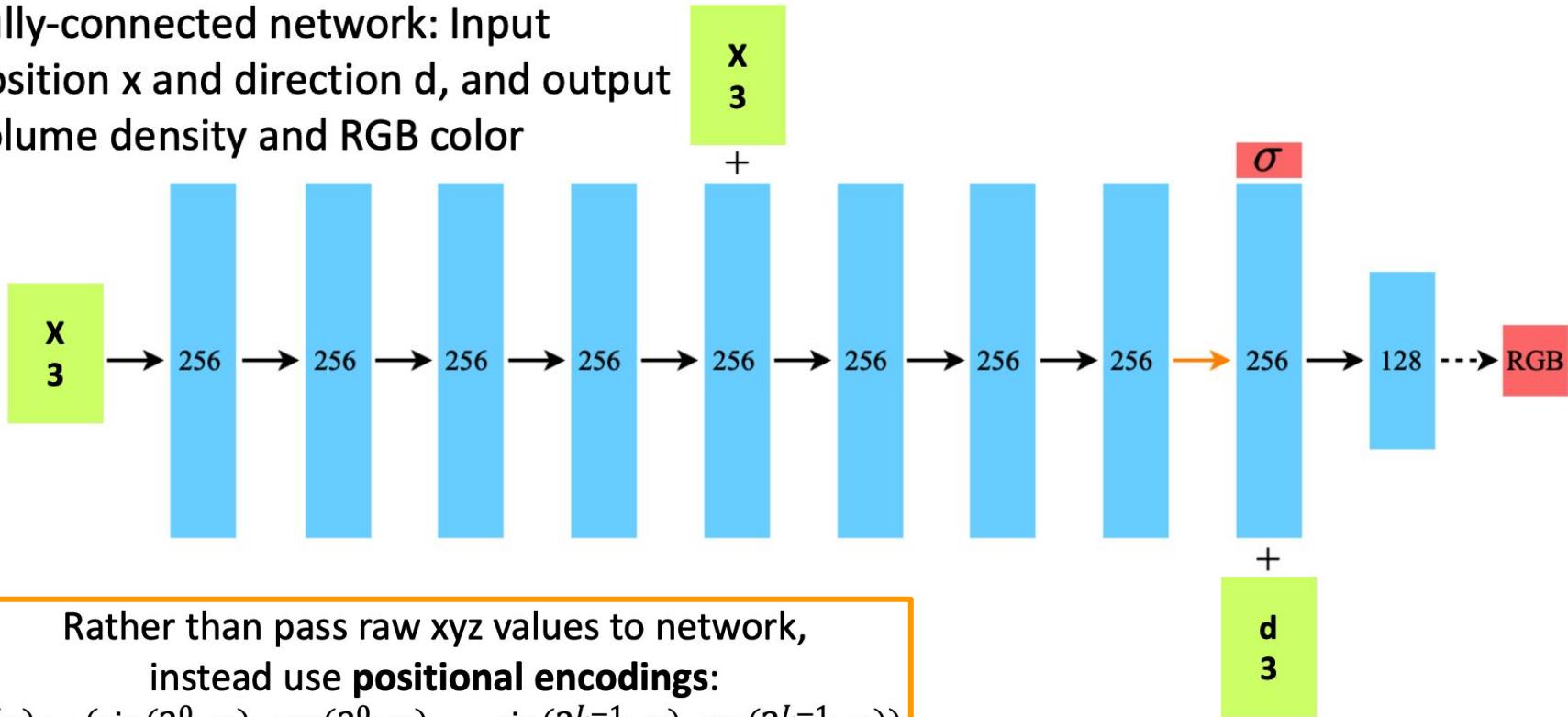
View-dependent
RGB color

Squared error

$$\mathcal{L} = \sum_{\mathbf{r} \in \mathcal{R}} \left[\left\| \hat{C}_c(\mathbf{r}) - C(\mathbf{r}) \right\|_2^2 + \left\| \hat{C}_f(\mathbf{r}) - C(\mathbf{r}) \right\|_2^2 \right]$$

NeRF: Network Architecture

Fully-connected network: Input position x and direction d , and output volume density and RGB color



NeRF: Network Architecture

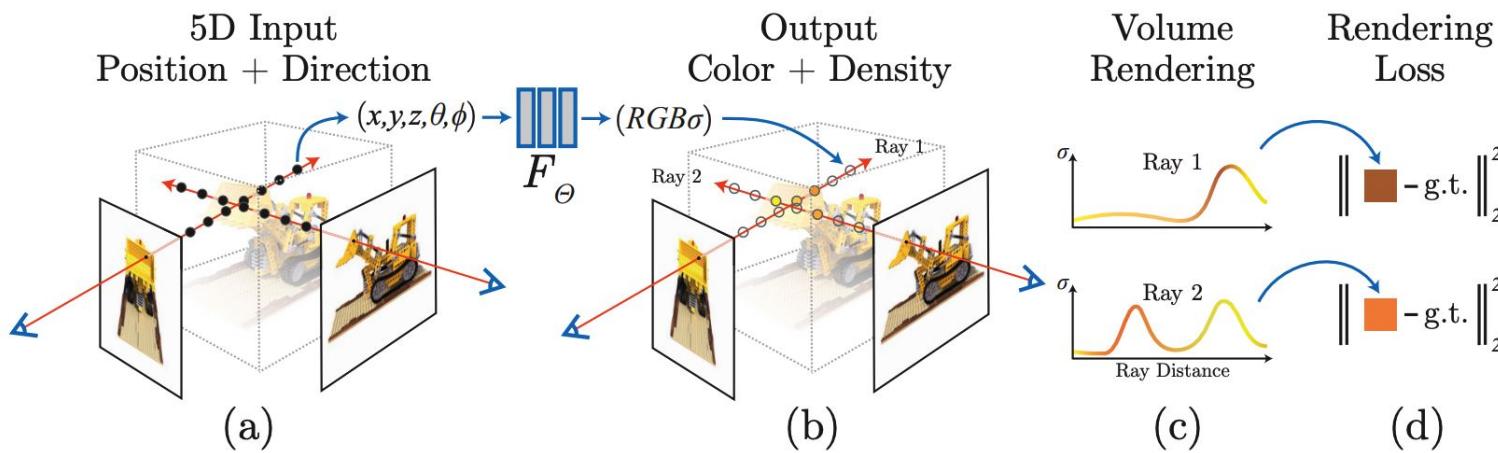


Fig. 2: An overview of our neural radiance field scene representation and differentiable rendering procedure. We synthesize images by sampling 5D coordinates (location and viewing direction) along camera rays (a), feeding those locations into an MLP to produce a color and volume density (b), and using volume rendering techniques to composite these values into an image (c). This rendering function is differentiable, so we can optimize our scene representation by minimizing the residual between synthesized and ground truth observed images (d).

NeRF

<https://www.matthewtancik.com/nerf>



NeRF (2020)

Main Problem: Very slow!

Also, not generalizable

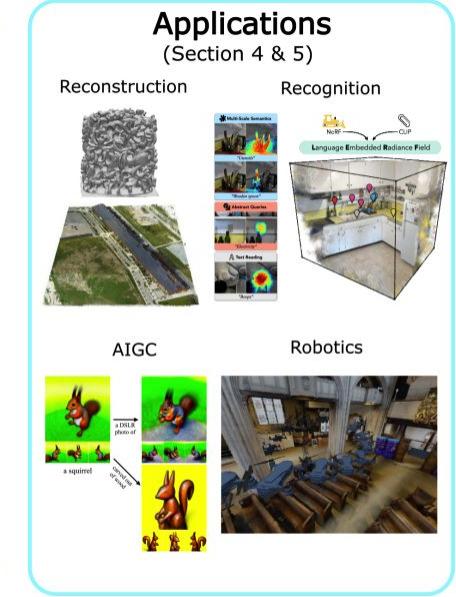
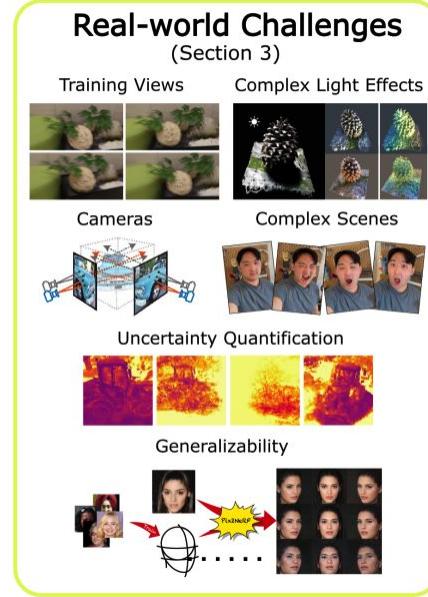
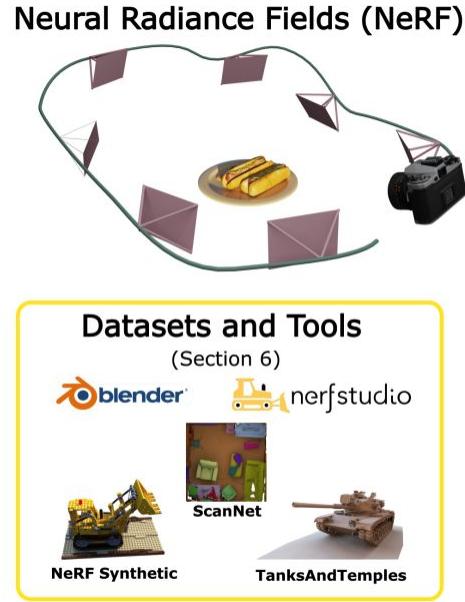
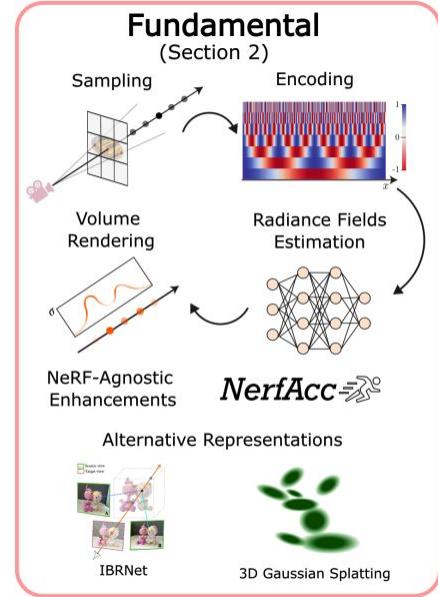
Training: 1-2 days on a V100 GPU, for just a single scene!

Inference: Sampling an image from a trained model:
 $(256 \times 256 \text{ pixels}) \times (224 \text{ samples per pixel})$
= 14.6M forward passes through MLP

Tons of follow-up work!

Cited by 11521

More on NeRF



Neural radiance field for the real world: a survey <https://arxiv.org/pdf/2501.13104v1.pdf>

Radiance Fields <https://radiancfields.com/>

Nerf and 3DGS at CVPR 2024 <https://github.com/Yubel426/NeRF-3DGS-at-CVPR-2024>

RS-Nerf (rolling shutter), ECCV 2024 <https://arxiv.org/pdf/2407.10267.pdf>

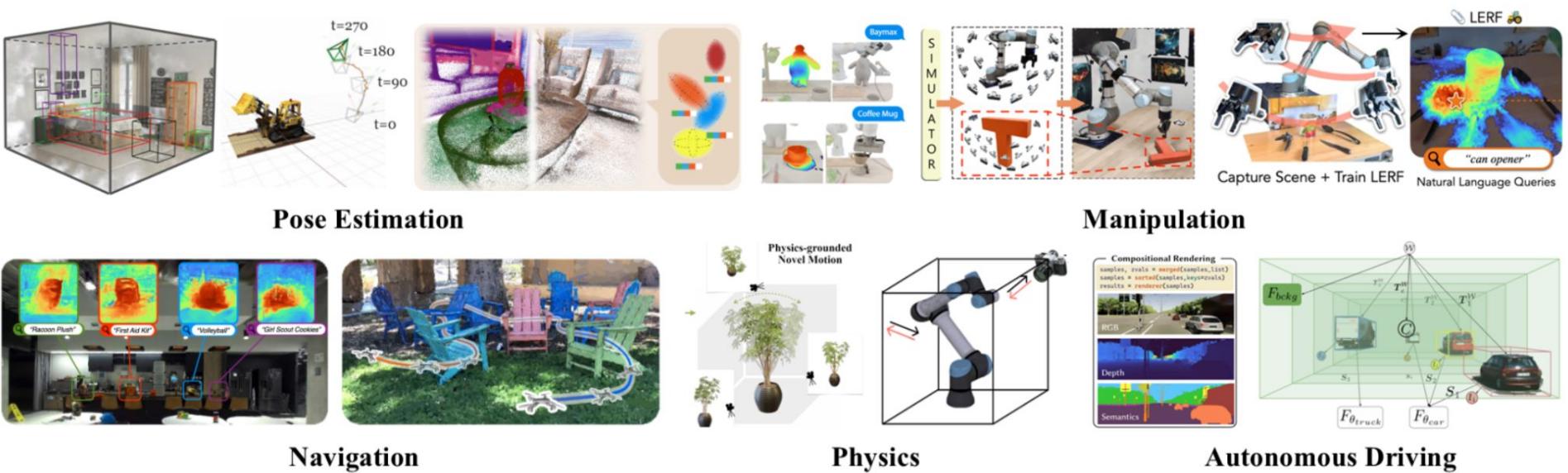
More on NeRF

Mip-NeRF360 <https://jonbarron.info/mipnerf360/> (unbounded scenes)

Block-NeRF <https://waymo.com/research/block-nerf/>

NeRF in Robotics <https://robonerf.github.io/survey/index.html>

LeRF <https://www.lerf.io/>

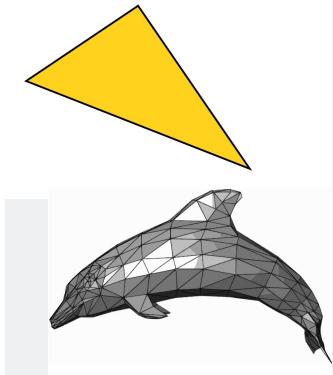
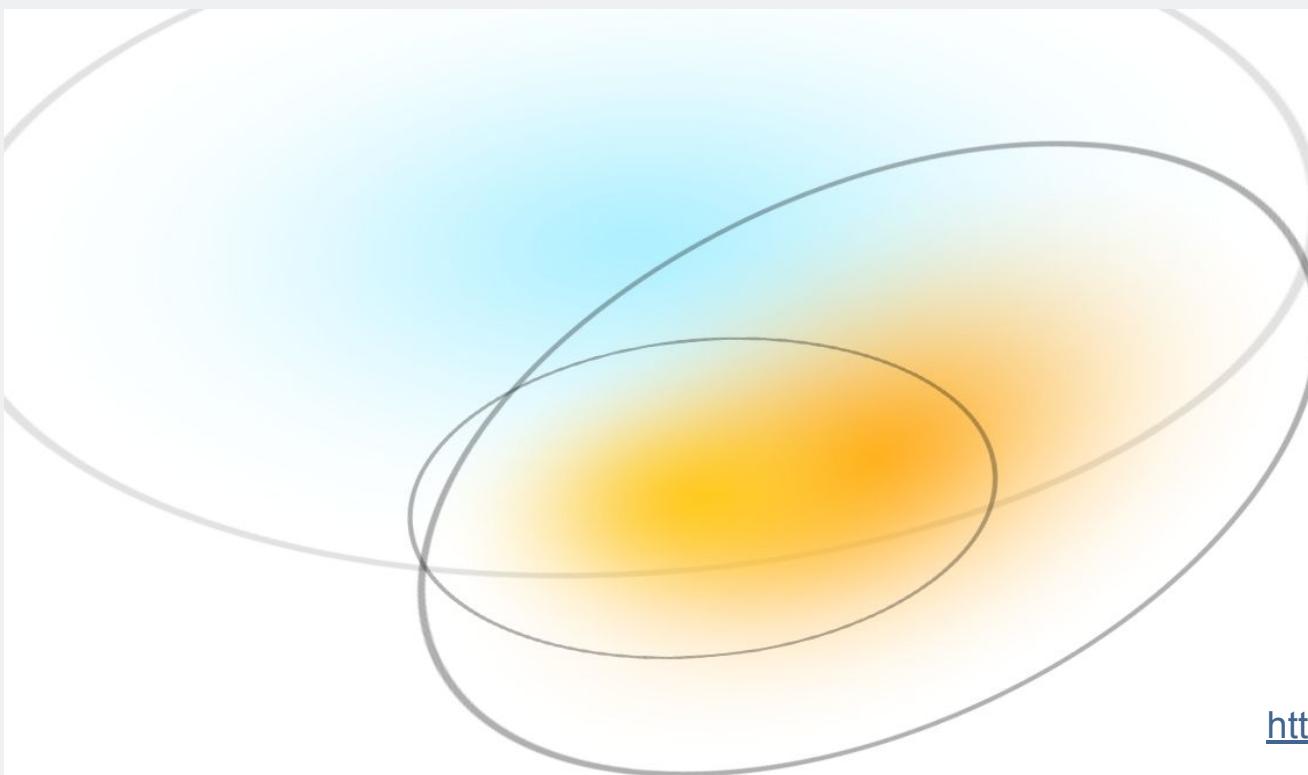


Gaussian Splatting

<https://arxiv.org/pdf/2308.04079>

Gaussian Splatting (2023)

Recall: Triangle Meshes



<https://arxiv.org/pdf/2308.04079>

Gaussian Splatting



Review: SfM

Structure from Motion

Input: a set of images
Output: image projections for each point
(camera parameters) + structure/scene
("geometrically verified images")



Figure 1. Result of Rome with 21K registered out of 75K images.

Gaussian Splatting

Input:

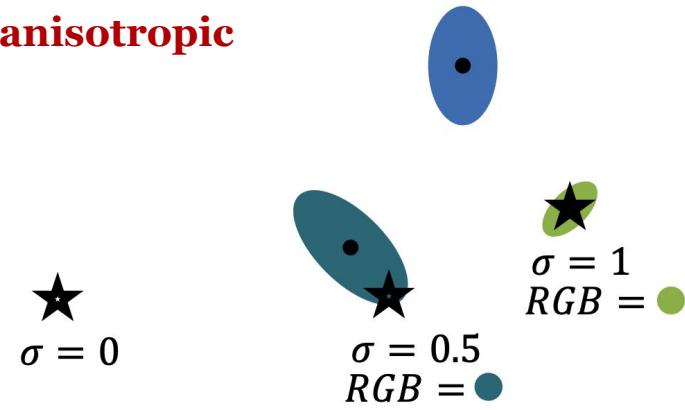
- a set of images of static scenes
- camera parameters (calibrated by SfM)

→ 3D Gaussians (**sparse**)

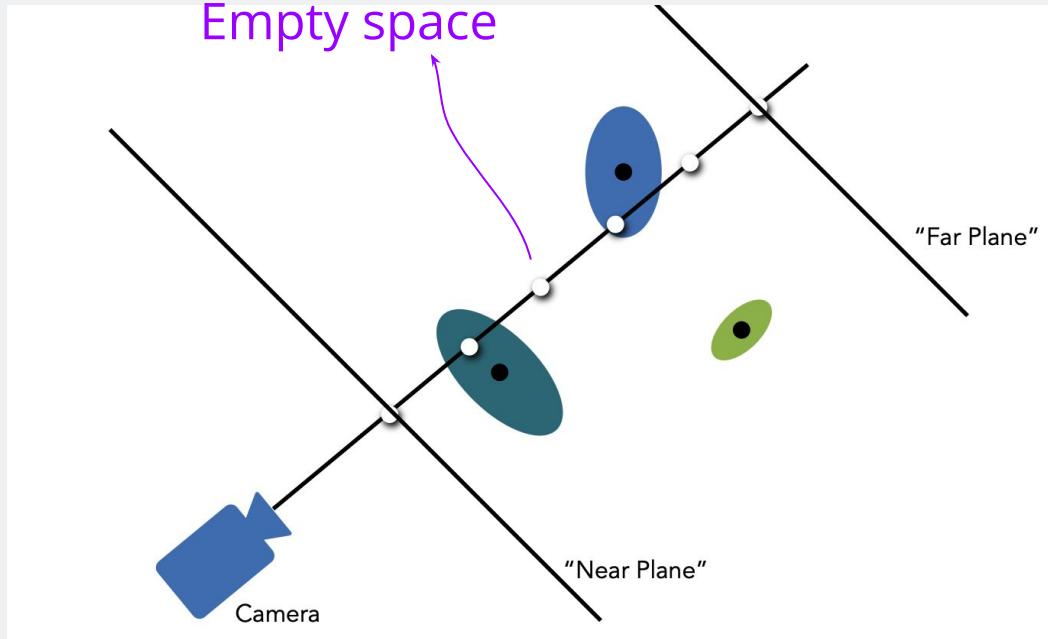
a position (mean), covariance matrix and opacity α

spherical harmonic (SH) coefficients (color)

anisotropic



Gaussian Splatting

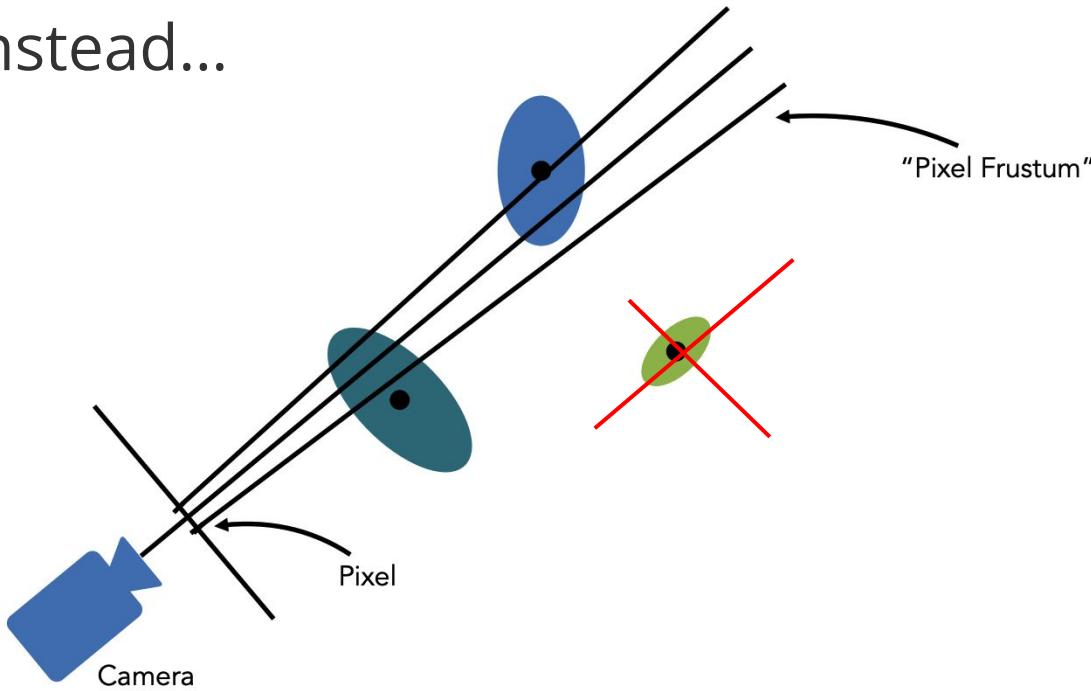


Similar to NeRF,
Same volume
rendering integral

BUT - let's NOT
sample a lot of
empty spaces...

Gaussian Splatting

Instead...



- Splitting the screen into 16×16 tiles
- “Cull” Gaussians against the **view frustum** and each tile
- Only keep Gaussians with a 99% confidence interval intersecting the view frustum

Additionally, adaptive Gaussian densification

<https://www.scenerepresentations.org/courses/2023/fall/inverse-graphics/>

https://web.stanford.edu/class/cs231a/lectures/lecture17_gaussian_splatting.pdf

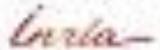


This video contains a voice-over.

3D Gaussian Splatting for Real-Time Radiance Field Rendering

SIGGRAPH 2023
(ACM Transactions on Graphics)

Bernhard Kerl^{*}



UNIVERSITÉ
CÔTE D'AZUR

Georgios Karanikas^{*}



UNIVERSITÉ
CÔTE D'AZUR

Thomas Leinweber



George Drettakis



UNIVERSITÉ
CÔTE D'AZUR

* Denotes equal contribution

More on Splat

3DGS tutorial <https://3dgstutorial.github.io/>

3D reconstruction <https://arxiv.org/pdf/2503.17316.pdf> (Pow3R, Dust3R)

4D Gaussian Splatting (March 31, 2025)

<https://arxiv.org/pdf/2503.22159.pdf>

Segmentation (March 28, 2025)

<https://vulab-ai.github.io/Segment-then-Splat/>

<https://arxiv.org/pdf/2503.22268.pdf> Segment any Motion (CVPR 2025)

Generative Models

Discriminative vs Generative Models

Discriminative Model:

Learn a probability distribution $p(y|x)$

Generative Model:

Learn a probability distribution $p(x)$

Conditional Generative Model: Learn $p(x|y)$

Data: x



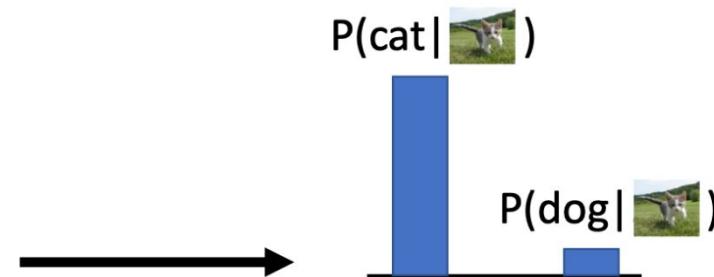
Label: y

Cat

Discriminative vs Generative Models

Discriminative Model:
Learn a probability distribution $p(y|x)$

classifiers



Generative Model:
Learn a probability distribution $p(x)$

Conditional Generative Model: Learn $p(x|y)$

Density Function
 $p(x)$ assigns a positive number to each possible x ; higher numbers mean x is more likely

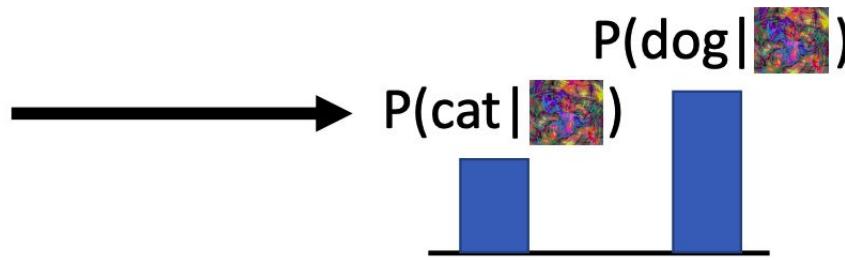
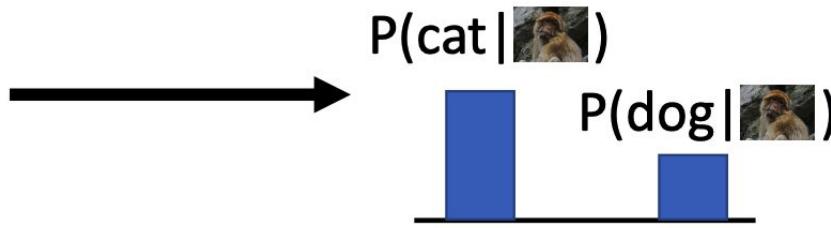
Density functions are **normalized**:

$$\int_X p(x)dx = 1$$

Different values of x **compete** for density

Discriminative vs Generative Models

Note about discriminative model:



Usually,

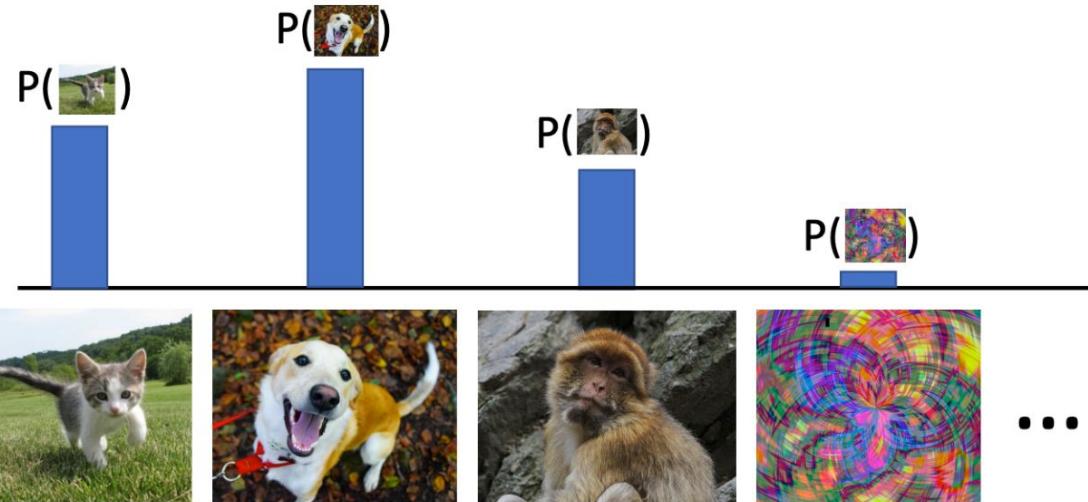
- No way for the model to handle unreasonable inputs;
- it must give label distributions for all images
- No “competition” between images

Discriminative vs Generative Models

Discriminative Model:
Learn a probability distribution $p(y|x)$

Generative Model:
Learn a probability distribution $p(x)$

Conditional Generative Model: Learn $p(x|y)$



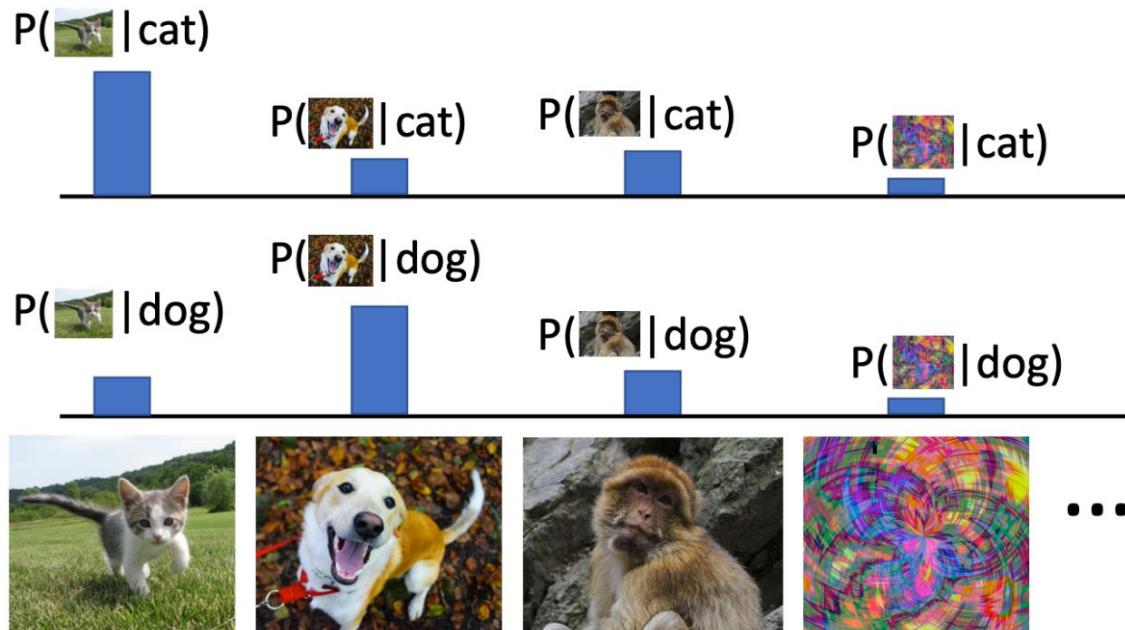
Generative model: All possible images compete with each other for probability mass

Discriminative vs Generative Models

Discriminative Model:
Learn a probability distribution $p(y|x)$

Generative Model:
Learn a probability distribution $p(x)$

Conditional Generative Model: Learn $p(x|y)$



Conditional Generative Model: Each possible label induces a competition among all images

Recall: Bayes' Rule

$$P(x | y) = \frac{P(y | x)}{P(y)} P(x)$$

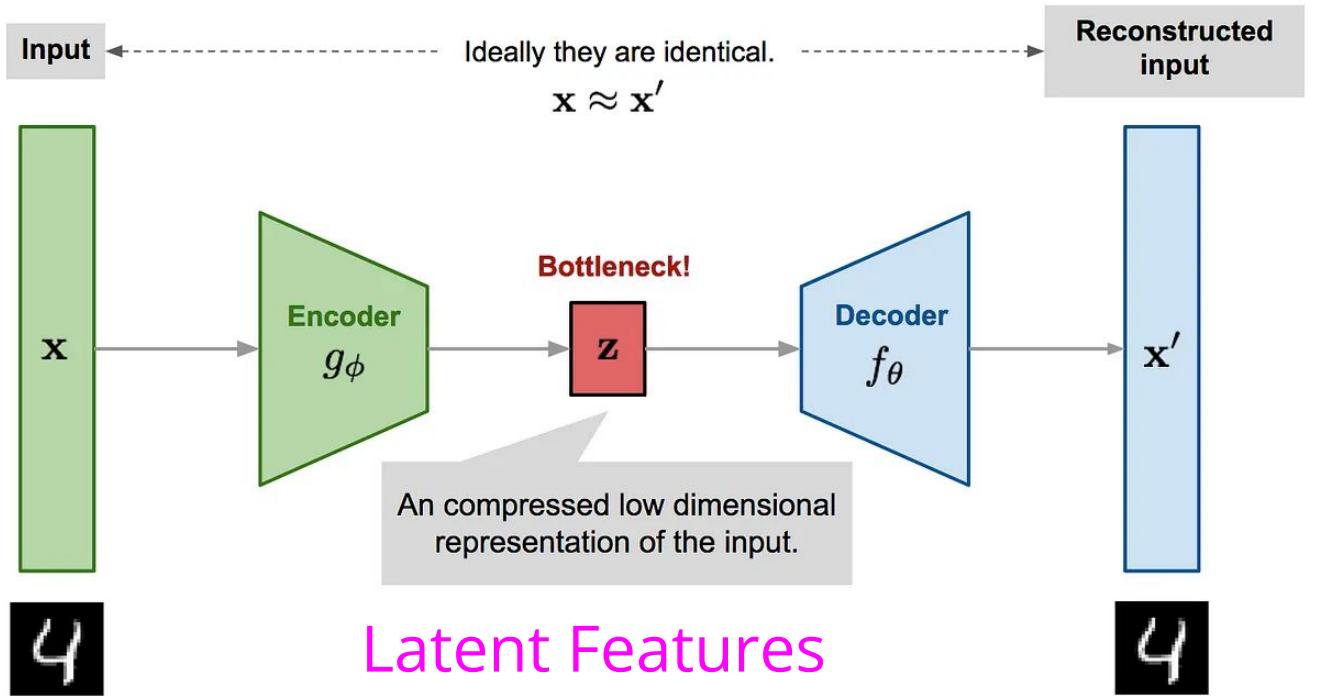
Conditional Generative Model Discriminative Model (Unconditional) Generative Model

Prior over labels

We can build a conditional generative model from other components!

AutoEncoders

Learn from data **without** labels



<https://lilianweng.github.io/posts/2018-08-12-vae/>

AutoEncoders

Learn from data **without** labels

Does not use any
labels! Just raw data!

Latent
Features
(low-dimension)

Input data

Reconstructed
input data

Loss Function

$$\|\hat{x} - x\|_2^2$$

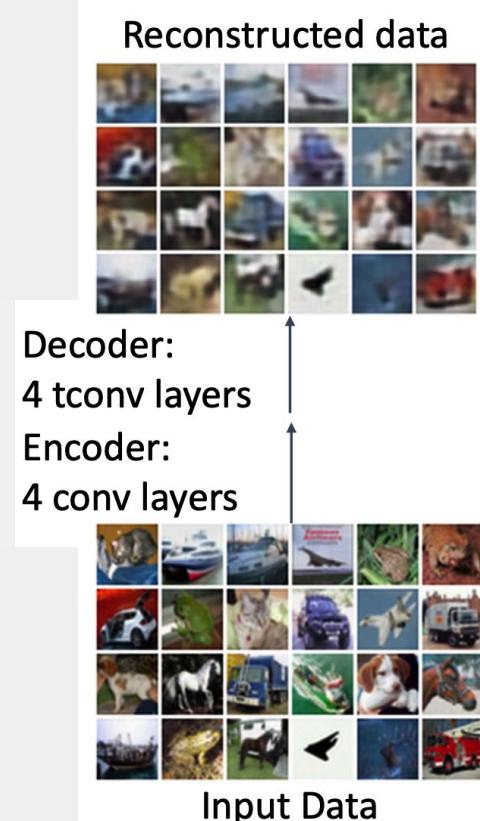
\hat{x}

Decoder

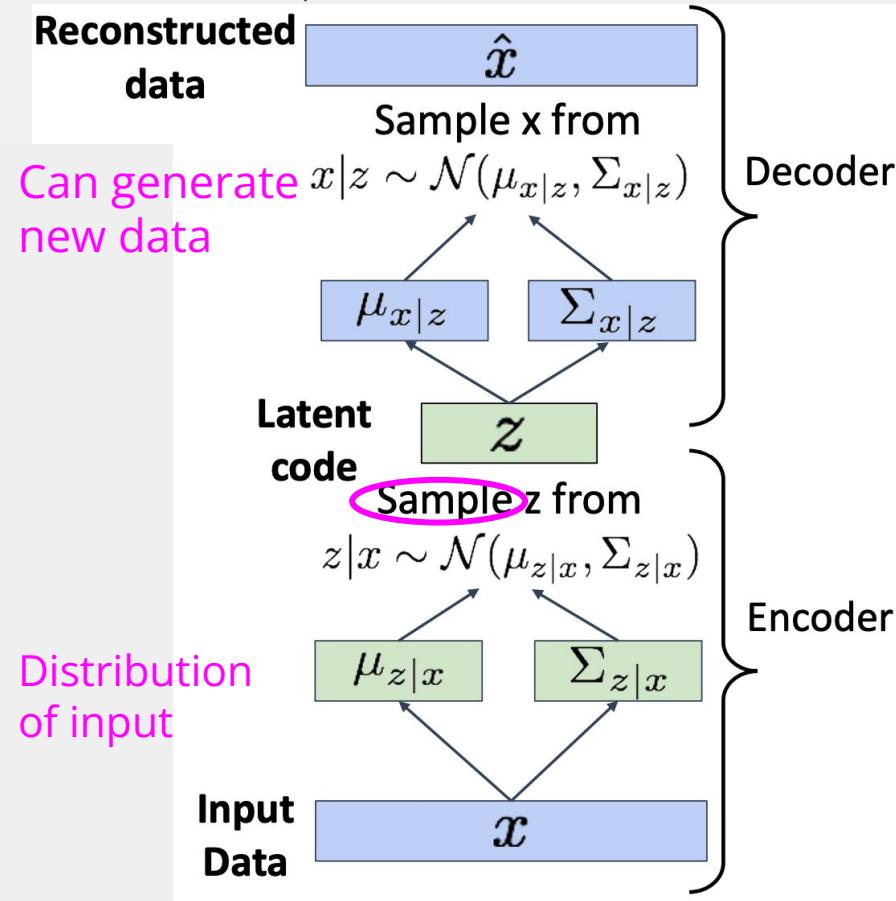
z

Encoder

x



VAE (Variational Autoencoders)



Train by maximizing the variational lower bound

$$E_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)] - D_{KL} (q_\phi(z|x), p(z))$$

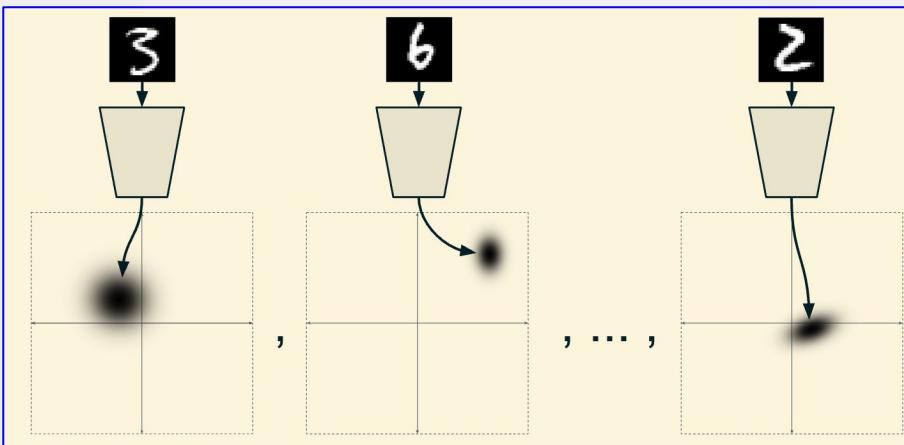
1. Run input data through **encoder** to get a distribution over latent codes
2. **Encoder output should match the prior $p(z)$!**
3. Sample code z from encoder output
4. Run sampled code through **decoder** to get a distribution over data samples
5. **Original input data should be likely under the distribution output from (4)!**
6. Can sample a reconstruction from (4)

<https://ijdykeman.github.io/ml/2016/12/21/cvae.html>

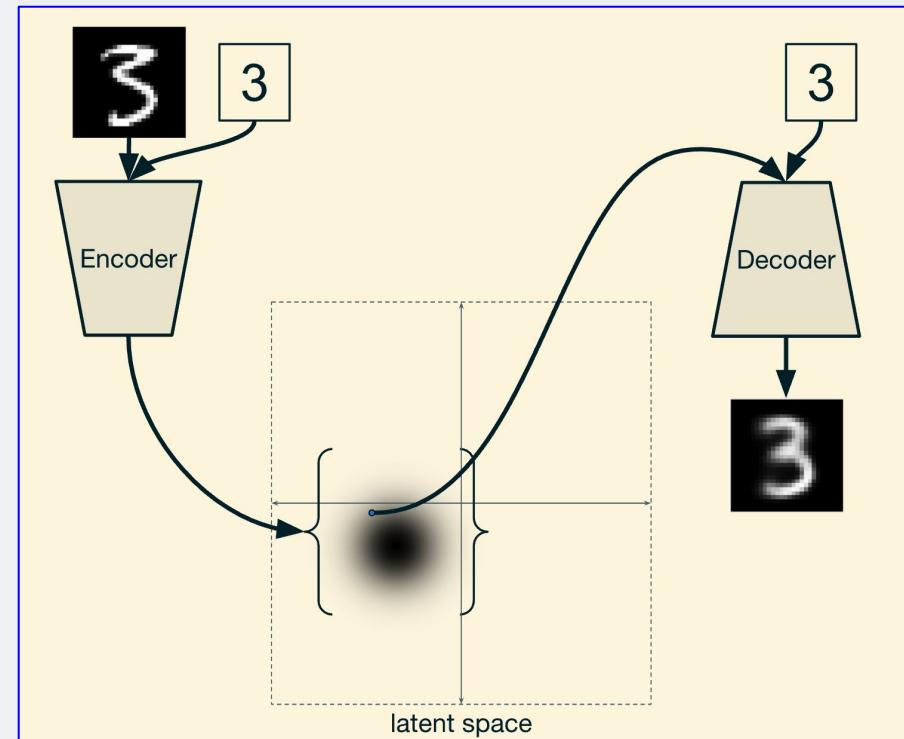
CVAE - labeled inputs

CVAE: Conditional VAE

VAE



CVAE - labeled inputs



GANs (Generative Adversarial Networks)

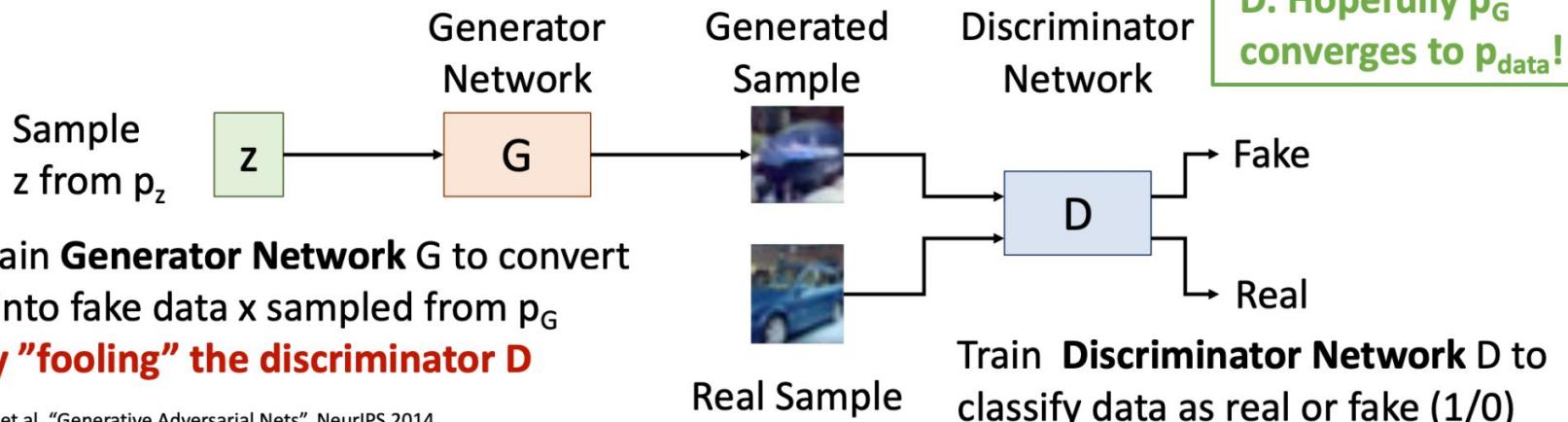
Setup: Assume we have data x_i drawn from distribution $p_{\text{data}}(x)$. Want to sample from p_{data} .

Idea: Introduce a latent variable z with simple prior $p(z)$.

Sample $z \sim p(z)$ and pass to a **Generator Network** $x = G(z)$

Then x is a sample from the **Generator distribution** p_G . Want $p_G = p_{\text{data}}$!

Jointly train G and D. Hopefully p_G converges to p_{data} !

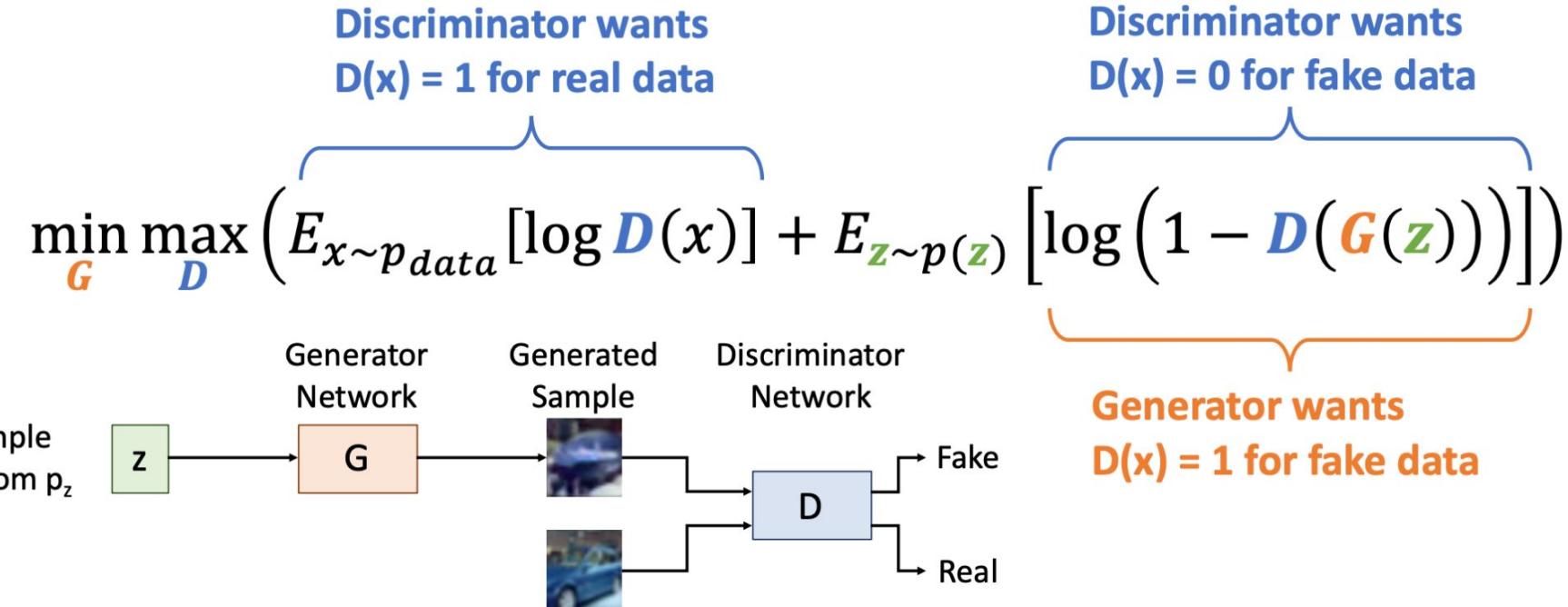


Goodfellow et al, "Generative Adversarial Nets", NeurIPS 2014

https://papers.nips.cc/paper_files/paper/2014/file/f033ed80deb0234979a61f95710dbe25-Paper.pdf

GANs (Generative Adversarial Networks)

Jointly train generator G and discriminator D with a **minimax game**



GANs (Generative Adversarial Networks)



Stable Diffusion

Noise → Data

v3.5, released Oct. 22, 2024

<https://stability.ai/news/introducing-stable-diffusion-3-5>

Stable diffusion v3 research paper

<https://arxiv.org/pdf/2403.03206.pdf>

<https://stability.ai/news/stable-diffusion-3-research-paper>

“outperforms state-of-the-art text-to-image generation systems such as DALL·E 3, Midjourney v6, and Ideogram v1 in typography and prompt adherence, based on human preference evaluations.”

Stable Diffusion



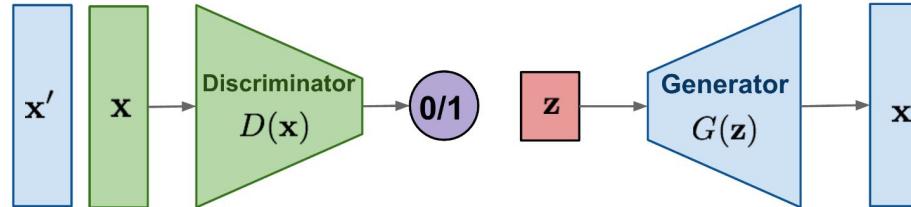
Noise → Data

Also, text prompts → images

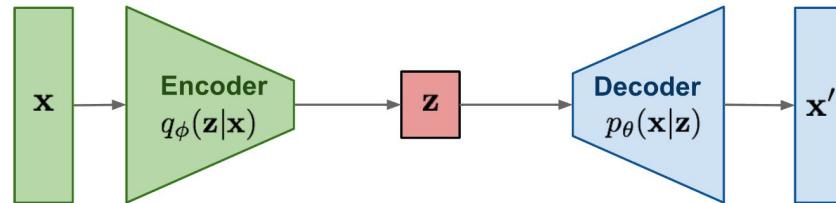
<https://scholar.harvard.edu/binxuw/classes/machine-learning-scratch/materials/stable-diffusion-scratch>

Diffusion Model

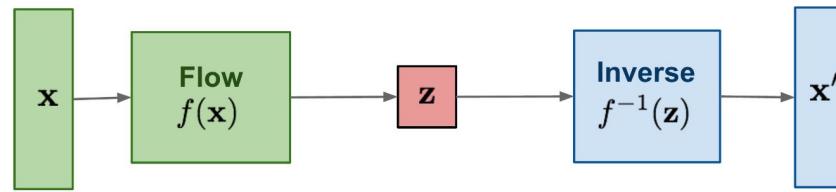
GAN: Adversarial training



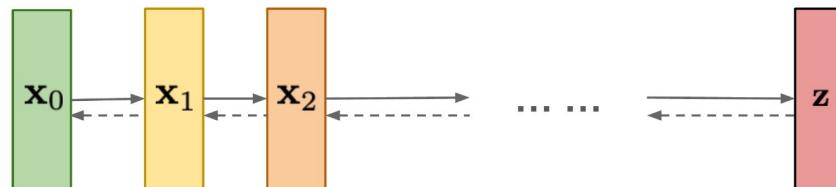
VAE: maximize variational lower bound



Flow-based models:
Invertible transform of distributions

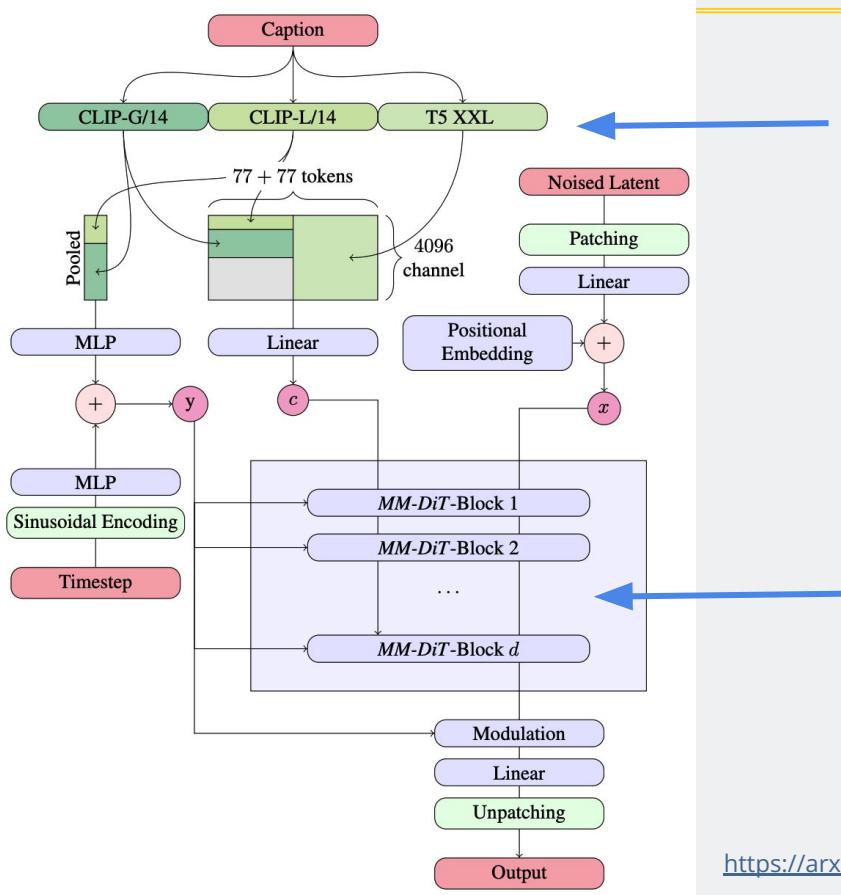


Diffusion models:
Gradually add Gaussian noise and then reverse



[/lilianweng.github.io/posts/2021-07-11-diffusion-models/](https://lilianweng.github.io/posts/2021-07-11-diffusion-models/)

Stable Diffusion



- Text/word representation
e.g., CLIP (more on this)
- Understanding prompts

Transformer attention modules
Text \longleftrightarrow Image

<https://arxiv.org/pdf/2403.03206>

Diffusion Model



1. Forward diffusion: **noising** $x_0 \rightarrow x_1 \rightarrow \dots x_T$

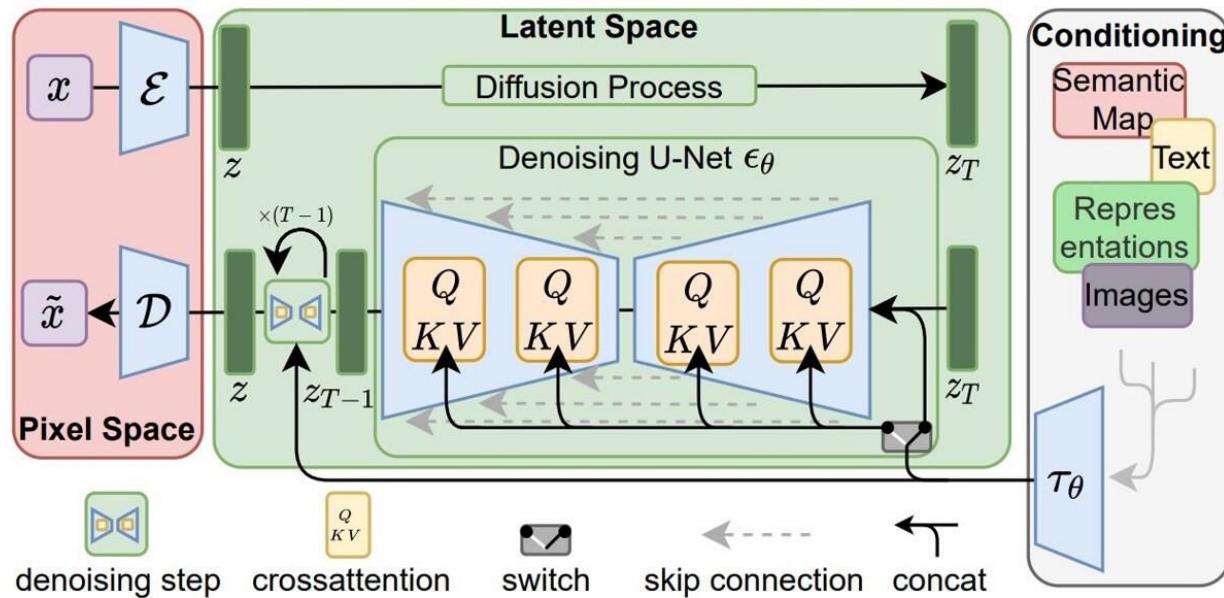
Take a data distribution $x_0 \sim p(x)$

Turn it into noise by $x_T \sim N(0, \sigma^2 I)$

2. Reverse diffusion: **denoising**

Sample from the noise distribution, reverse the diffusion process to generate data

Diffusion Model



- Principle of Diffusion models (sampling, learning)
- Diffusion for Images – e.g., UNet architecture
- Understanding prompts – Word as vectors, e.g, CLIP
- Let words modulate diffusion – Conditional Diffusion, Cross Attention
- Diffusion in latent space – AutoEncoderKL
- Training on **Massive Dataset**

SORA

<https://openai.com/sora/>



ROBOTICS

“Video generation models as world simulators”

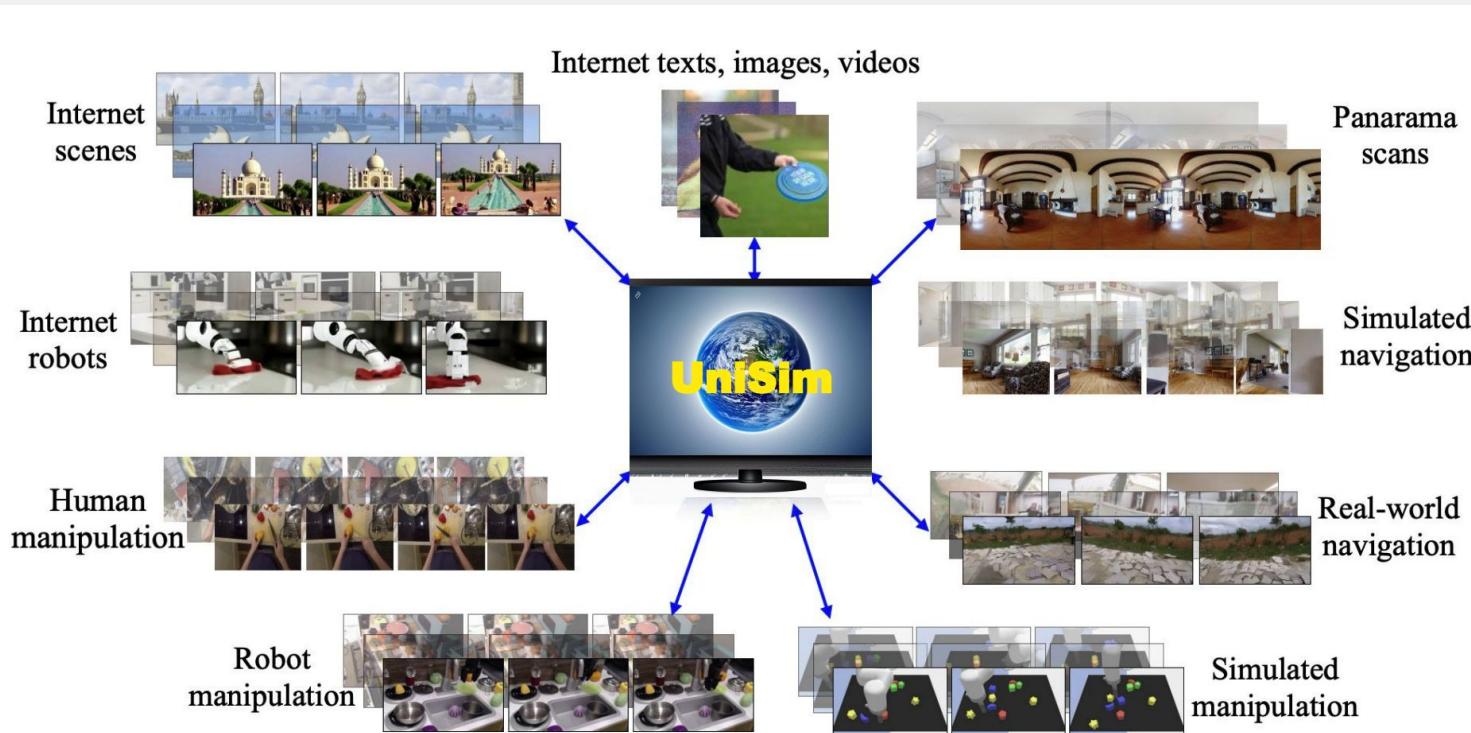
<https://openai.com/index/video-generation-models-as-world-simulators/>

Not official technical report from openAI <https://arxiv.org/pdf/2402.17177>

Uni-Sim

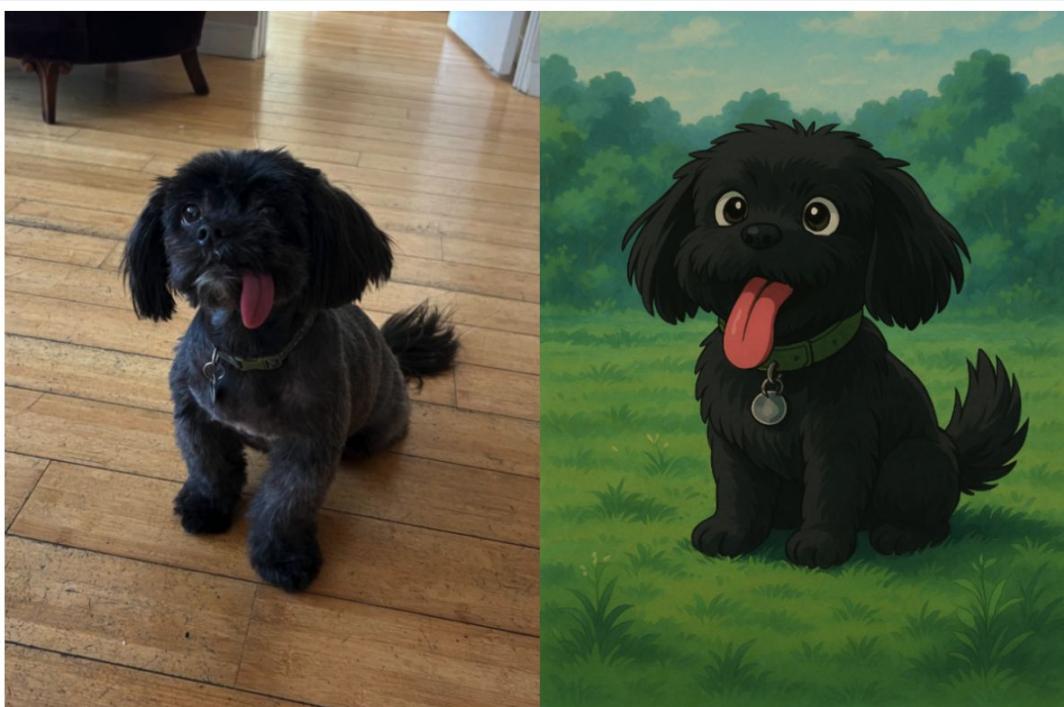
<https://universal-simulator.github.io/unisim/> (ICLR 2024)

<https://arxiv.org/pdf/2310.06114>



“Action-Rich”
Simulations

Image Generation



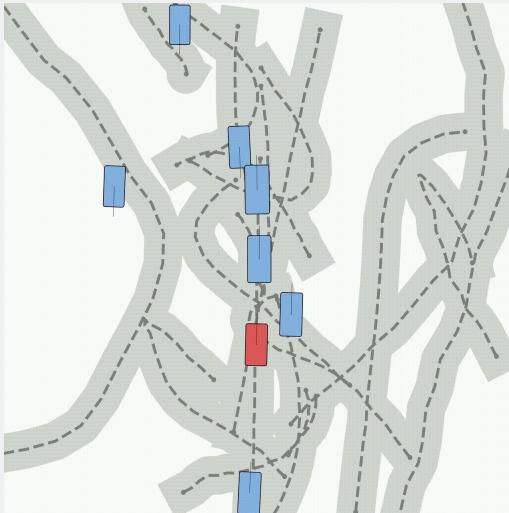
A REAL DOG (LEFT) AND AN AI-GENERATED IMAGE OF A DOG CHATGPT MADE IN STUDIO GHIBLI STYLE (RIGHT).

IMAGE CREDITS:MAXWELL ZEFF/OPENAI

More on Diffusion

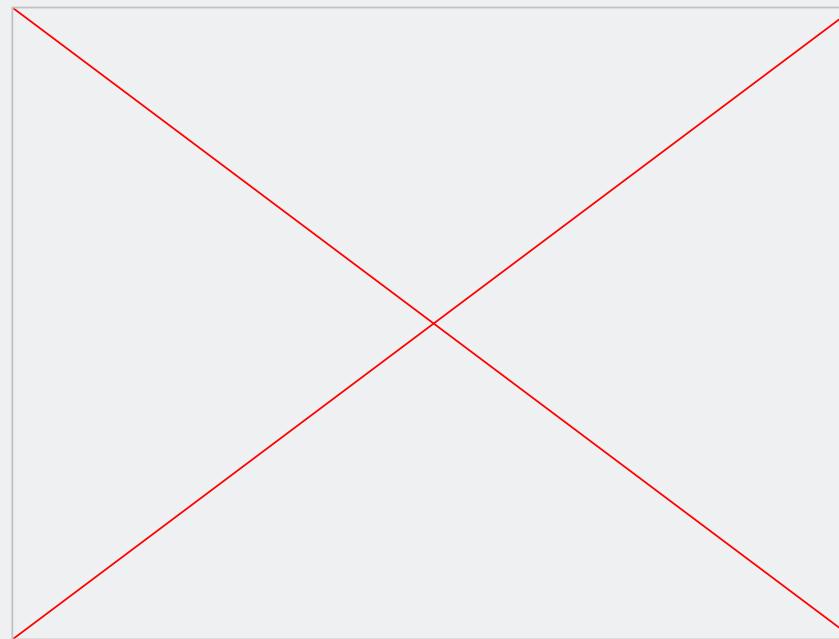
Scenario Dreamer

<https://princeton-computational-imaging.github.io/scenario-dreamer/> (CVPR 2025)



Stable virtual camera

<https://stable-virtual-camera.github.io/>



Robotic Piano playing

<https://taco-group.github.io/PANDORA/>

Summary

