# Vision AI Demystified

Ayesha Asif, Developer Marketing  |  GTC March 2024

# Agenda

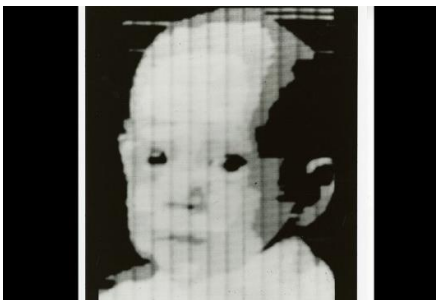- The History of Computer Vision

- From CNNs to Vision Transformers

- Generative AI: GANs and Diffusion Models

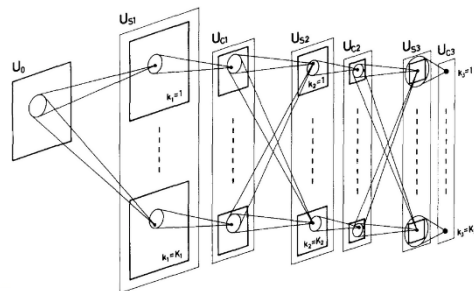- 3D Reconstruction with NeRFs and Gaussian Splatting

- Emerging Computer Vision Frontiers

nVIDIA

# A Brief History of Computer Vision



**1957**
**The first digitally scanned photos**
Russell Kirsch

**1963**
**Machine Perception of Three-Dimensional Solids**
Lawrence Roberts

**1970**
**Neocognitron**
Kunihiko Fukushima

**1989**
**LeNet-5**
Yann LeCun

**2001**
**The Viola-Jones algorithm**
Paul Viola and Michael Jones

**2009**
**ImageNet Dataset**
Fei-Fei Li

**2012-2024**
**Deep Learning**
Vision Transformers
Diffusion Models
Generative AI
3D Reconstruction

# Convolutional Neural Networks



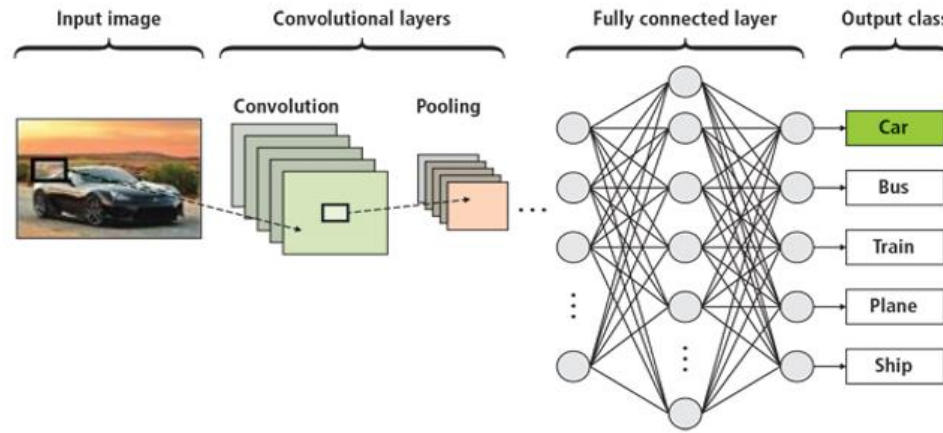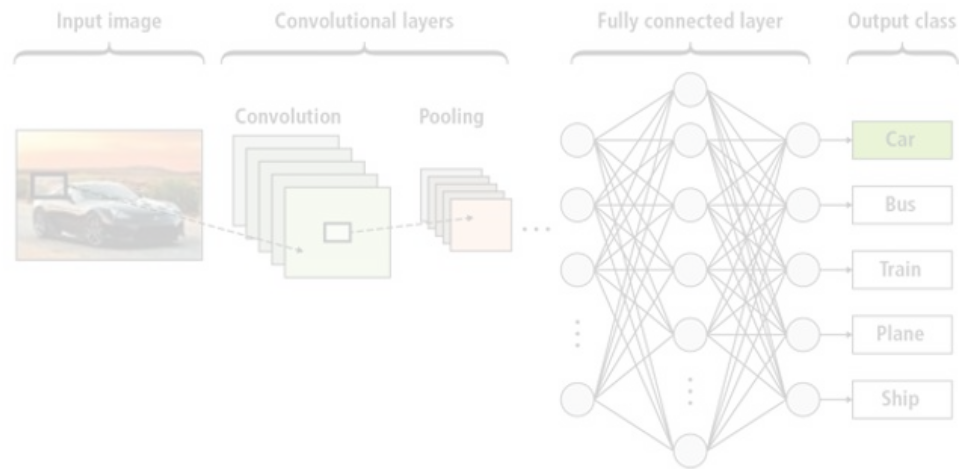## CNN Architecture

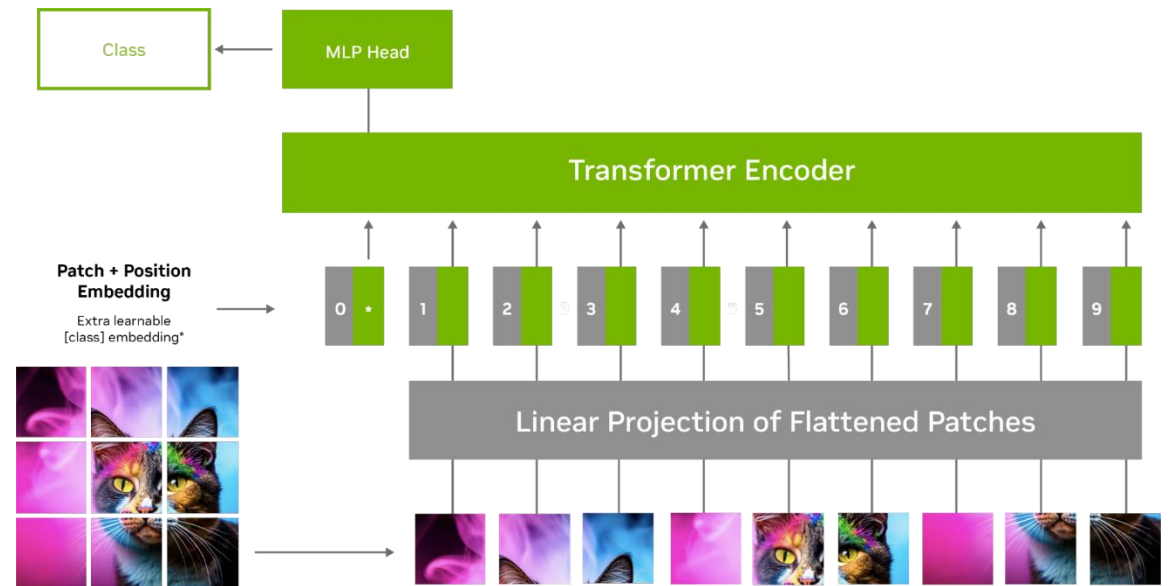Hierarchical feature extraction and classification

# Convolutional Neural Networks and Transformers

## What are Vision Transformers (ViTs)?



**CNN Architecture**

Hierarchical feature extraction and classification

**Vision Transformer Architecture**

Attention-based patch embedding and classification

# Comparison of CNN and ViT Methods

# Vision Transformers

loss:
$- p_2 \log p_1$

DINO

- Self-supervised learning method

- Same architecture for teacher and student networks

- Student network – matches the output distribution of a teacher network

- Teacher network – learns from the student parameters through exponential moving average (EMA)

# Vision Transformers

Segment Anything Model (SAM)



SAM

- SAM was trained on an unprecedented dataset of 11 million images and 1.1 billion segmentation masks

- Strong zero-shot performance

- Promptable segmentation

# Vision Transformers

## Contrastive Language-Image Pretraining (CLIP)



CLIP

- Trained on 400 million (image, text) pairs

- Jointly trains an image and text encoder to predict the correct pairings of a (image, text) examples

- Projected in a shared embedding space to understand the text/image relationships

- Applications – image search, image captioning

# Getting Started With Vision Transformers



## DINO Notebook

#Installing the TAO launcher
```
!pip3 install nvidia-pyindex
!pip3 install nvidia-tao
```

#Pull pre-trained model from NGC
```
!ngc registry model download-version
nvidia/tao/pretrained_dino_nvimagenet:fan_small_hybri
d_nvimagenet --dest $LOCAL_PROJECT_DIR/dino/
```

#Run Inference of a pre-trained model
```
tao model dino inference -e /path/to/spec.yaml -r
/path/to/results/
inference.checkpoint=/path/to/model.pth
```
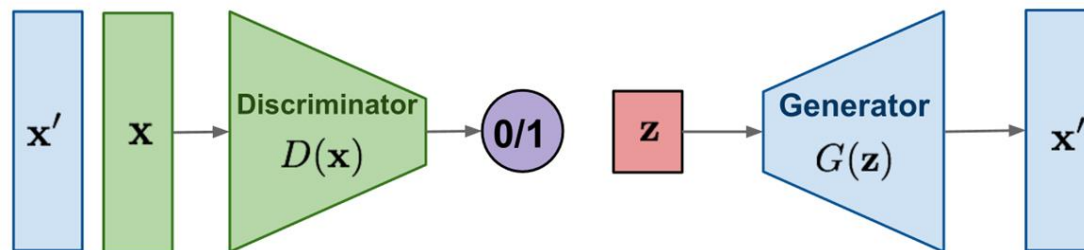
Finetuning on Custom Datasets:
- Tao model dino train
- Tao model dino evaluate

# Generative Models

# Generative Models

## Generative Adversarial Networks (GANs)



**GAN Architecture**

Adversarial training

- The generator creates synthetic data, and the discriminator tries to distinguish the generated data from real data

- Training instability and mode collapse issues

# Generative Models

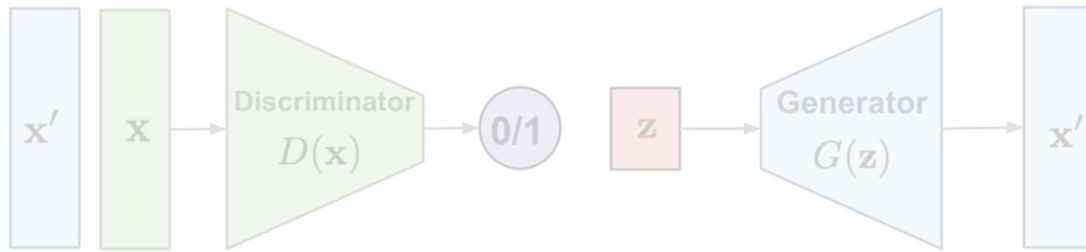## From GANs to Diffusion Models



GAN Architecture

Adversarial training

- The generator creates synthetic data, and the discriminator tries to distinguish the generated data from real data

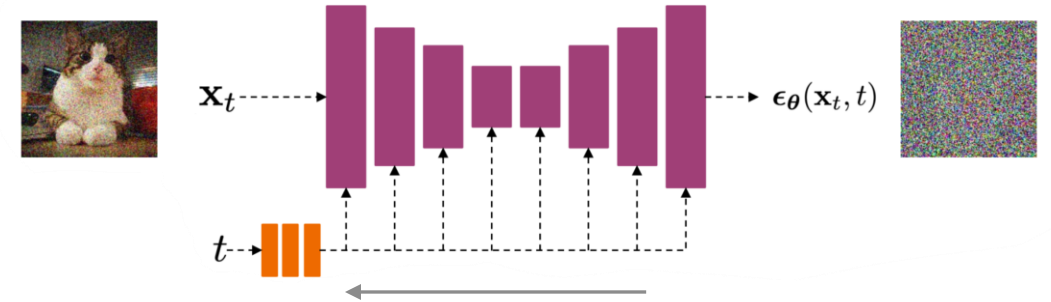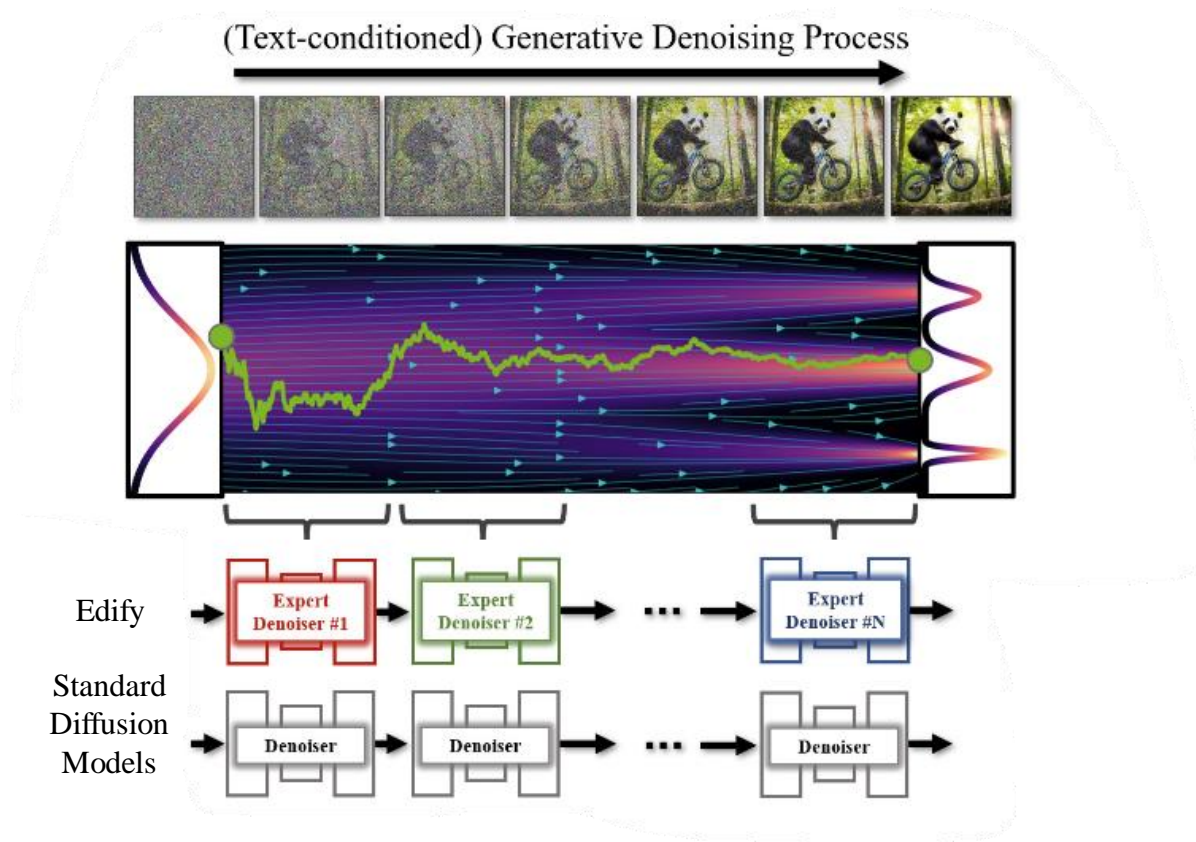- Training instability and mode collapse issues

Diffusion Model Architecture

Gradual addition and removal of noise

- Learn to reverse a diffusion process by iteratively removing noise from data

- Can generate high-quality samples across various domains

# Generative Models

NVIDIA Edify



(Text-conditioned) Generative Denoising Process

Edify

Standard Diffusion Models



Credit : Getty Images AI Generator

Wildlife photography of wolf, beautiful eyes, golden hour

- Instead of using a single denoiser, Edify trains an ensemble of expert denoising networks
- Edify also combines text embeddings from both T5 and CLIP models instead of using just one text encoder

# Generative AI by **getty**images

Built using NVIDIA Picasso

- State-of-the-art NVIDIA Edify model architecture for 4K generative photography

- Commercially safe—trained exclusively on Getty Images licensed data; uncapped indemnification

- Access via Web or iStock.com

- Advanced visual editing APIs for adding subjects, expanding images or replacing specific element

NVIDIA

nature photography of a rock arch, a mountain lake with a forest and fog in the background, overcast, majestic

# Getting Started With Diffusion Models



A photo of a Shiba Inu dog with a backpack riding a bike

```python
import requests

# Call model from NVIDIA NGC
invoke_url = "https://api.nvcf.nvidia.com/v2/nvcf/pexec/functions/89848fb8-549f-41bb-88cb-95d6597044a4"
fetch_url_format = "https://api.nvcf.nvidia.com/v2/nvcf/pexec/status/"

# Provide API key here
headers = {
"Authorization": "Bearer $API_KEY_REQUIRED_IF_EXECUTING_OUTSIDE_NGC",
"Accept": "application/json",
}

# Provide prompts and parameters
payload = {
"prompt": "A photo of a Shiba Inu dog with a backpack riding a bike",
"negative_prompt": "beach",
"inference_steps": 25
}

# re-use connections
session = requests.Session()

response = session.post(invoke_url, headers=headers, json=payload)

while response.status_code == 202:
        response = session.get(fetch_url_format + response.headers.get("NVCF-REQID"), headers=headers)

response.raise_for_status()
response_body = response.json()
print(response_body)
```

catalog.ngc.nvidia.com/orgs/nvidia/teams/ai-foundation/models/sdxl

# Getting Started With Diffusion Models



A photo of a Shiba Inu dog with
a backpack riding a bike

```python
import requests

# Call model from NVIDIA NGC
invoke_url = "https://api.nvcf.nvidia.com/v2/nvcf/pexec/functions/89848fb8-549f-41bb-88cb-95d6597044a4"
fetch_url_format = "https://api.nvcf.nvidia.com/v2/nvcf/pexec/status/"

# Provide API key here
headers = {
"Authorization": "Bearer $API_KEY_REQUIRED_IF_EXECUTING_OUTSIDE_NGC",
"Accept": "application/json",
}

# Provide prompts and parameters
payload = {
"prompt": "A photo of a Shiba Inu dog with a backpack riding a bike",
"negative_prompt": "beach",
"inference_steps": 25
}

# re-use connections
session = requests.Session()

response = session.post(invoke_url, headers=headers, json=payload)

while response.status_code == 202:
        response = session.get(fetch_url_format + response.headers.get("NVCF-REQID"), headers=headers)

response.raise_for_status()
response_body = response.json()
print(response_body)
```

catalog.ngc.nvidia.com/orgs/nvidia/teams/ai-foundation/models/sdxl

NVIDIA

# Getting Started With Diffusion Models



A photo of a Shiba Inu dog with
a backpack riding a bike

```python
import requests

# Call model from NVIDIA NGC
invoke_url = "https://api.nvcf.nvidia.com/v2/nvcf/pexec/functions/89848fb8-549f-41bb-88cb-95d6597044a4"
fetch_url_format = "https://api.nvcf.nvidia.com/v2/nvcf/pexec/status/"

# Provide API key here
headers = {
"Authorization": "Bearer $API_KEY_REQUIRED_IF_EXECUTING_OUTSIDE_NGC",
"Accept": "application/json",
}

# Provide prompts and parameters
payload = {
"prompt": "A photo of a Shiba Inu dog with a backpack riding a bike",
"negative_prompt": "beach",
"inference_steps": 25
}

# re-use connections
session = requests.Session()

response = session.post(invoke_url, headers=headers, json=payload)

while response.status_code == 202:
        response = session.get(fetch_url_format + response.headers.get("NVCF-REQID"), headers=headers)

response.raise_for_status()
response_body = response.json()
print(response_body)
```

catalog.ngc.nvidia.com/orgs/nvidia/teams/ai-foundation/models/sdxl

# Getting Started With Diffusion Models

# 3D Reconstruction

# 3D Reconstruction

## Structure from Motion (SfM)



3D–Model

line of sight

image i

corresponding
feature points

image i+1

image i+2

image i+3

moving camera

Source: Theia Vision Library

- An imaging technique for estimating 3D structures from 2D image sequences

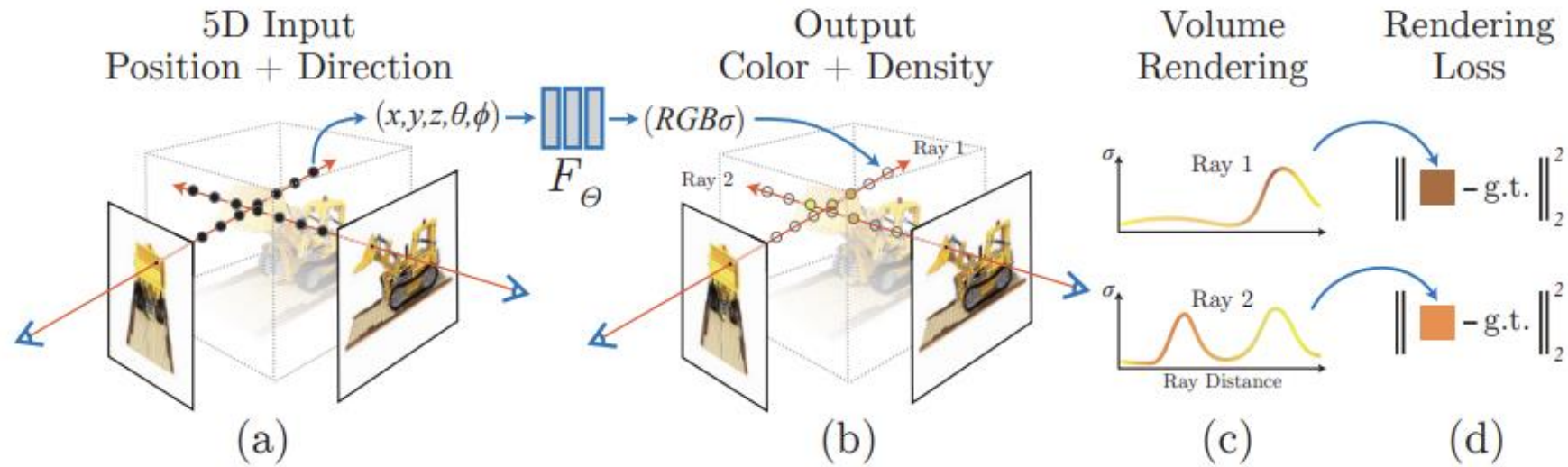- SfM tracks features (e.g. corner points, edges) across multiple images to find correspondences.

- This method struggles to capture fine details and complex geometry and struggles with certain materials

# 3D Reconstruction

## Neural Radiance Fields (NeRFs)



**NeRF architecture**

- NeRFs can generate novel views of complex 3D scenes, based on a partial set of 2D images

- Several decoupled modules - Sampling, encoding, neural network, and rendering

- Diverse applications spanning virtual/augmented reality, medical imaging, robotics, and autonomous navigation

# 3D Reconstruction

Neuralangelo

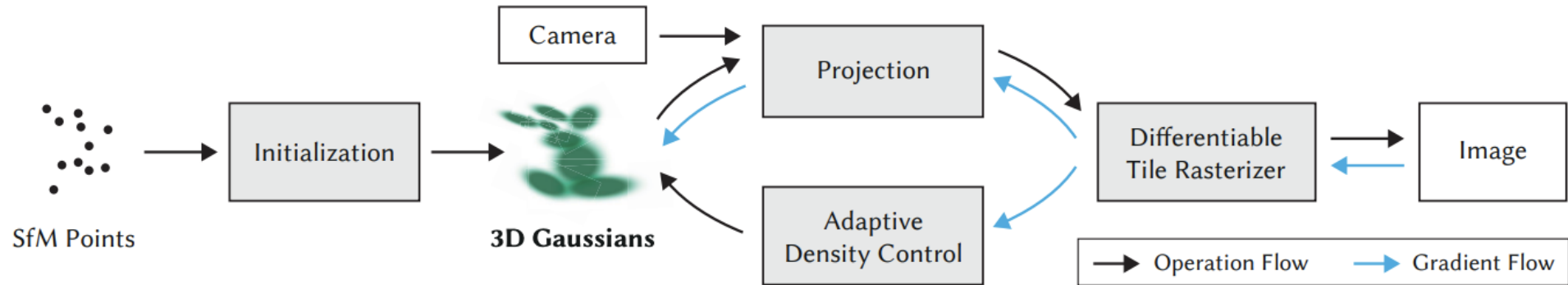# 3D Online Shopping Experiences Powered by NeRFs

RECON Labs Inc.



Take a video of the object

# 3D Reconstruction

## 3D Gaussian Splatting



3D Gaussian Splatting for Real-Time Radiance Field Rendering

# 3D Reconstruction

## 3D Gaussian Splatting



**Point Cloud Generation**
- 3D sparse point are generated using 3D points

**Gaussian Placement**
- Each point transformed into a 3D Gaussian function defined by position, spread, color and opacity

**Training & Optimization**
- Optimizes parameter of each gaussian points (position, spread, color, opacity)

**Rendering**
- Optimized points are projected on the screen based on camera's viewpoint
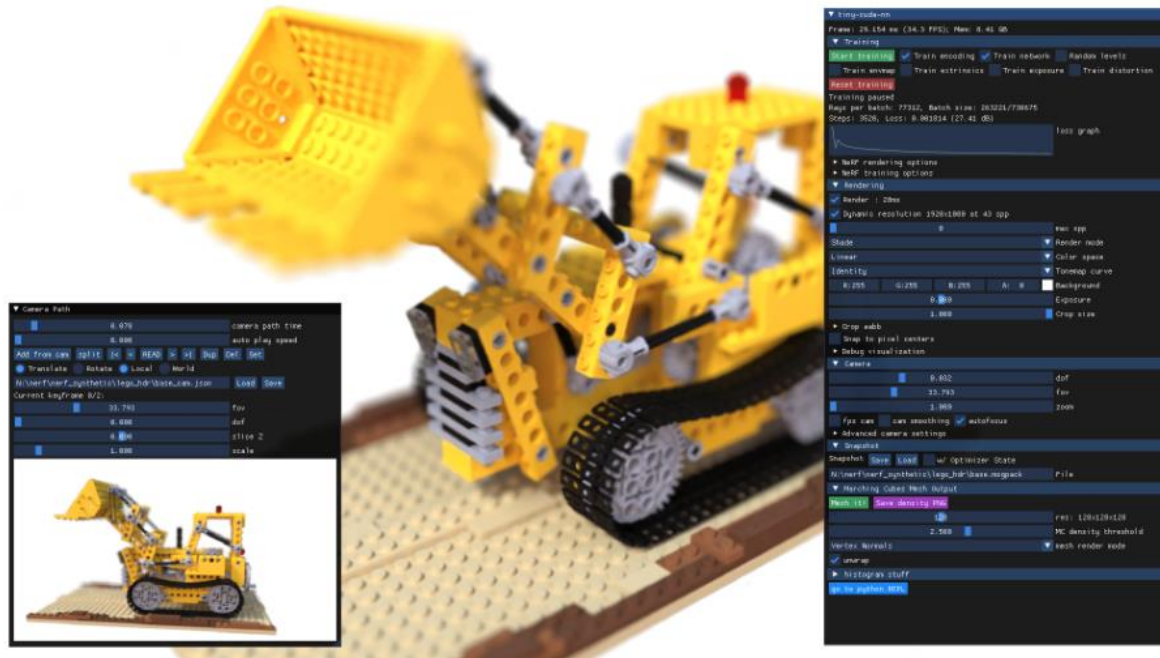
3D Gaussian Splatting for Real-Time Radiance Field Rendering

# 3D Reconstruction

## 3D Gaussian Splatting



Instant-NGP

Gaussian Splatter

3D Gaussian Splatting for Real-Time Radiance Field Rendering

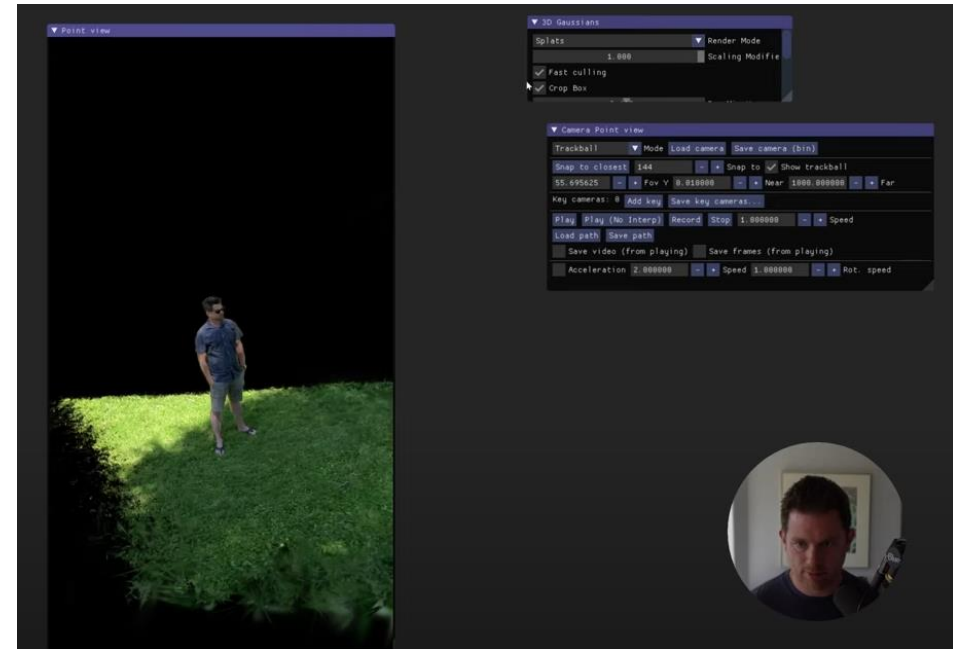# Getting Started With Instant-NGP & 3D Gaussian Splatting



Video Tutorial by Eric Hanes, NVIDIA

Step by Step Instructions in GitHub

Video Tutorial by The NeRF Guru

Step by Step Instructions in GitHub

# Emerging Computer Vision Frontiers

# Emerging Computer Vision Frontiers

Text-2-Image with Character Consistency

# Emerging Computer Vision Frontiers

Text-to-4D (3D in motion)

NVIDIA.

# Other Relevant GTC Sessions

- [S62724] Revolutionizing Vision AI: From 2D to 3D Worlds

- [S62818] The Visionaries: A Cross-Industry Exploration of Computer Vision

- [CWE63456] Build Accelerated Computer Vision Microservices With NVIDIA Libraries and SDKs

- [S62624] The Vision-AI Revolution powered by DeepStream