

Image and Video Generations

Lecture 0: Course Introduction

劉育綸

Yu-Lun (Alex) Liu

What do you know about image
and video generation models?



air head · Made by shy kids with Sora

OpenAI Sora



Luma Dream Machine Ray2



Black Forest Labs Flux



Adobe Firefly



AI-Generated Art Won a Prize



Jason Allen won the digital-art competition at the Colorado State Fair last year for his piece "Theatre D'opera Spatial" that he created using the AI software Midjourney. Recently, the US Copyright Office refused to grant him a copyright for his piece, writing, "We have decided that we cannot register this copyright claim because the deposit does not contain any human authorship.¹ He plans to appeal.

AI ASMR – Glass Fruits Cutting



Google DeepMind Veo 3



<https://deepmind.google/models/veo/>

Will Smith Eating Spaghetti test



2023

Will Smith Eating Spaghetti test



<https://www.instagram.com/reel/C3i5vAZvRS3/>

Will Smith Eating Spaghetti test



2023



2025

Google Nano Banana

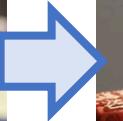


Image Edit Arena

Rank (UB) ↑	Model ↑	Score ↑	95% CI (±) ↑↓	Votes ↑	Organization ↑	License ↑↓
1	 gemini-2.5-flash-image-preview (nano-banana)	1362	±2	2,521,035	Google	Proprietary
2	 flux-1-kontext-max	1191	±3	357,196	Black Fores...	Proprietary
3	 flux-1-kontext-pro	1174	±2	2,015,530	Black Fores...	Proprietary
3	 gpt-image-1	1170	±3	1,026,399	OpenAI	Proprietary
5	 flux-1-kontext-dev	1152	±3	1,584,400	Black Fores...	Proprietary
6	 qwen-image-edit	1145	±2	1,585,904	Alibaba	Apache 2.0
6	 seededit-3.0	1142	±4	1,285,080	Bytedance	Proprietary
8	 gemini-2.0-flash-preview-image-generation	1093	±3	1,700,785	Google	Proprietary
9	 bagel	1044	±5	12,774	Bytedance	Apache 2.0
10	 steplx-edit	1017	±4	138,399	StepFun	Apache 2.0

Text-to-3D Generation



Vintage copper rotary telephone with intricate detailing.



Two-story brick house with red roof and fence.



Glowing orb on a stone pedestal.



Spherical robot with gold and silver design.



Bronze owl sculpture perched on a branch.



Futuristic robotic arm on a table.

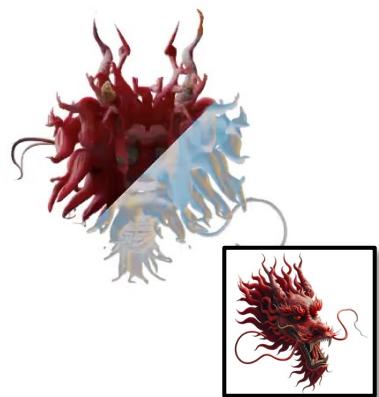
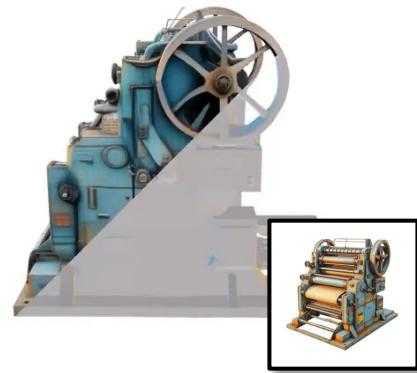


A rustic log cabin with a stone chimney and a wooden porch.



Blocky, orange and teal robot with articulated limbs.

Image-to-3D Generation



Text-to-4D Generation



A storm trooper walking forward and vacuuming, best quality, 4K, HD.



Beer pouring into a glas.



A bee fluttering its wings fast.



A cat singing, best quality, 4K, HD.



Mage in purple robe running forward, full body, portrait, game, unreal, 4K, HD.



Santa Claus carrying a big bag is walking forward, portrait, game, unreal, 4K, HD.



A dog wearing a Superhero outfit with red cape flying through the sky.



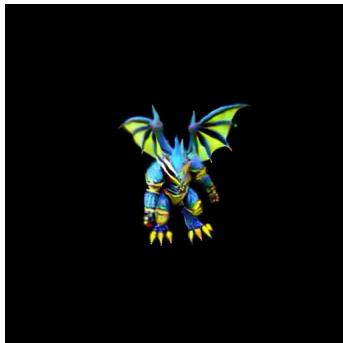
Flying dragon on fire.



An astronaut riding a horse, best quality, 4K, HD.



An ancient roman statue dancing, full body, portrait, game, unreal, 4K, HD.



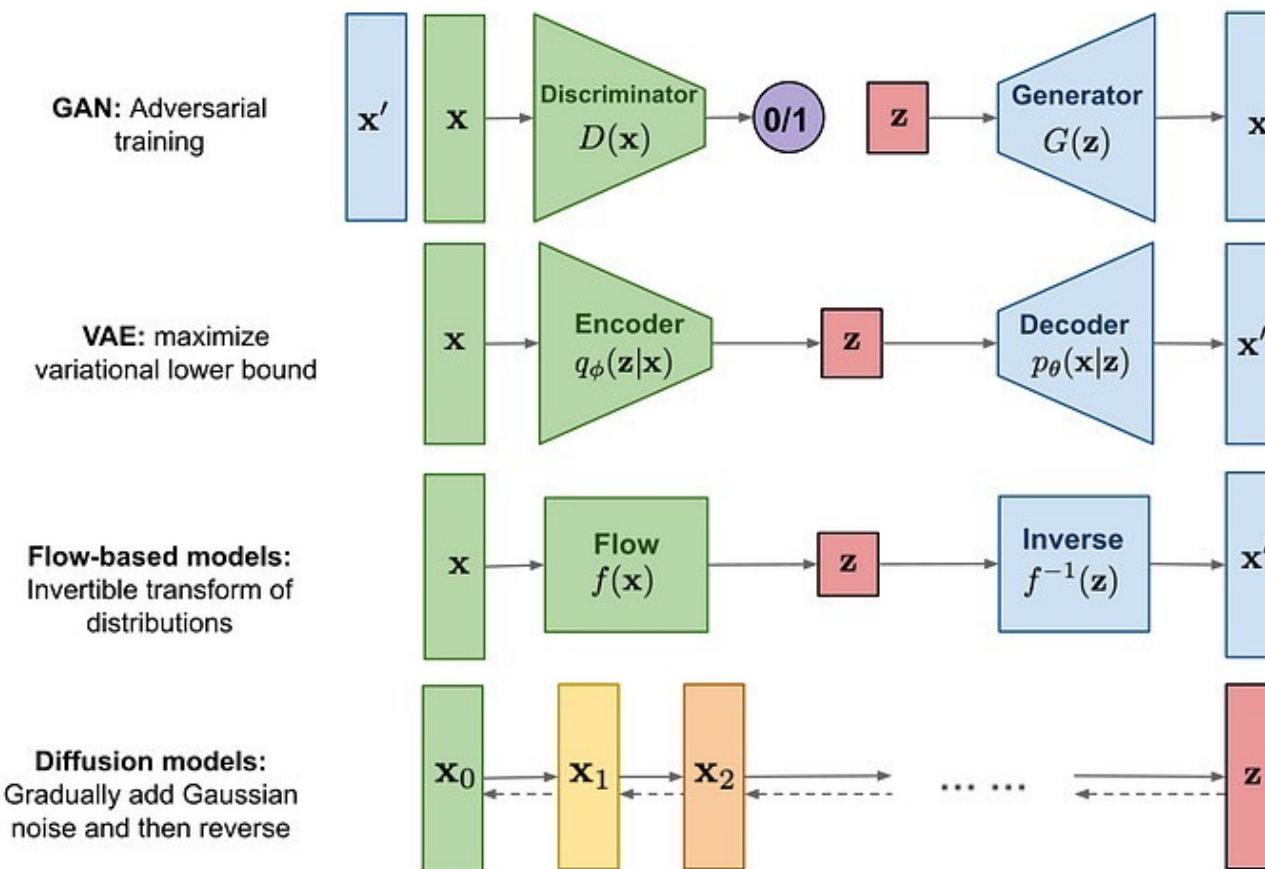
Dragon armor fluttering its wings fast.



A corgi is running fast.

A Brief Overview of the Course

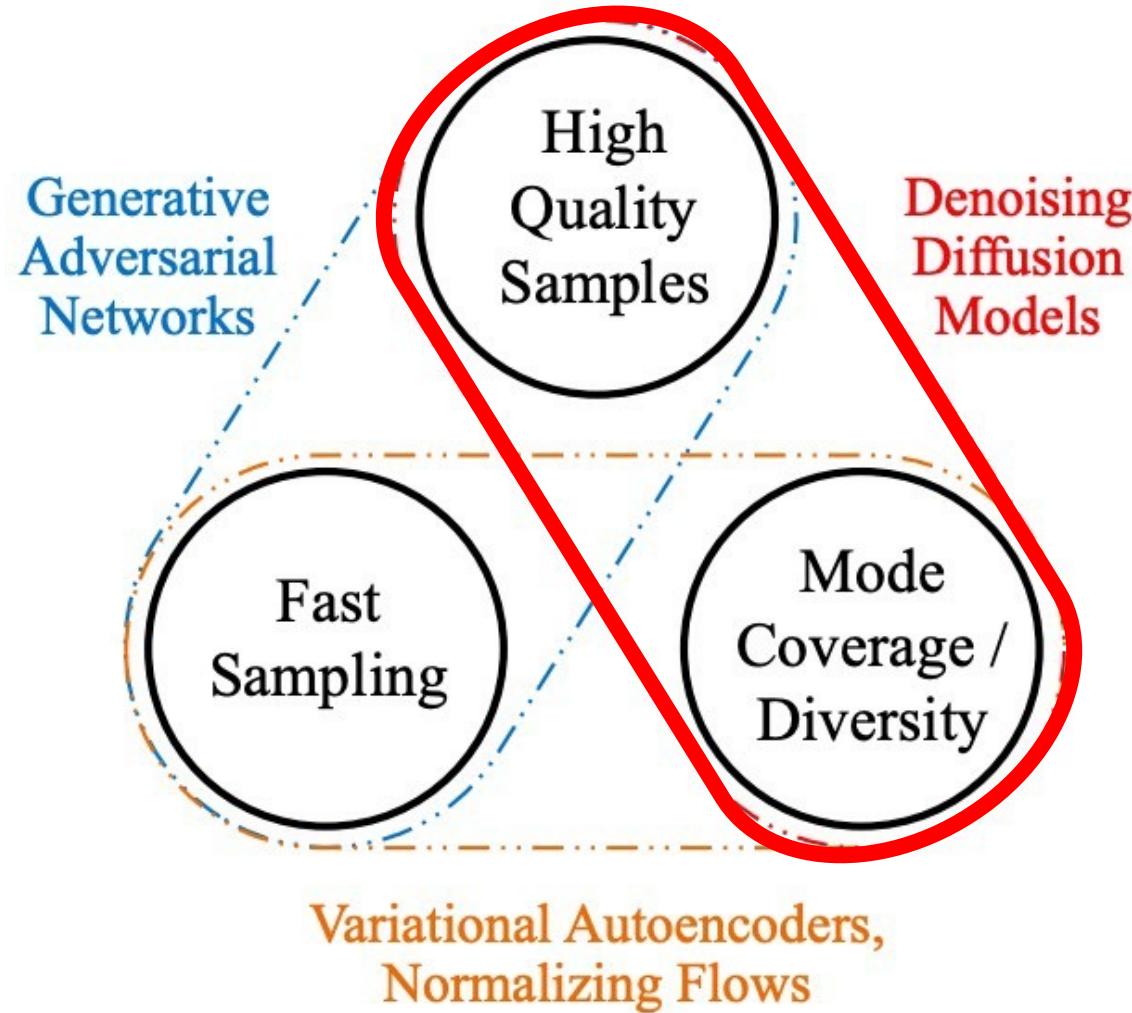
Generative Models



StyleGAN2



Generative Models – Comparison



Higher Quality & Diversity

Diffusion Models Beat GANs on Image Synthesis

Prafulla Dhariwal*

OpenAI

prafulla@openai.com

Alex Nichol*

OpenAI

alex@openai.com

Abstract

We show that diffusion models can achieve image sample quality superior to the current state-of-the-art generative models. We achieve this on unconditional image synthesis by finding a better architecture through a series of ablations. For conditional image synthesis, we further improve sample quality with classifier guidance: a simple, compute-efficient method for trading off diversity for fidelity using gradients from a classifier. We achieve an FID of 2.97 on ImageNet 128×128 , 4.59 on ImageNet 256×256 , and 7.72 on ImageNet 512×512 , and we match BigGAN-deep even with as few as 25 forward passes per sample, all while maintaining better coverage of the distribution. Finally, we find that classifier guidance combines well with upsampling diffusion models, further improving FID to 3.94 on ImageNet 256×256 and 3.85 on ImageNet 512×512 . We release our code at <https://github.com/openai/guided-diffusion>.

Diffusion Models

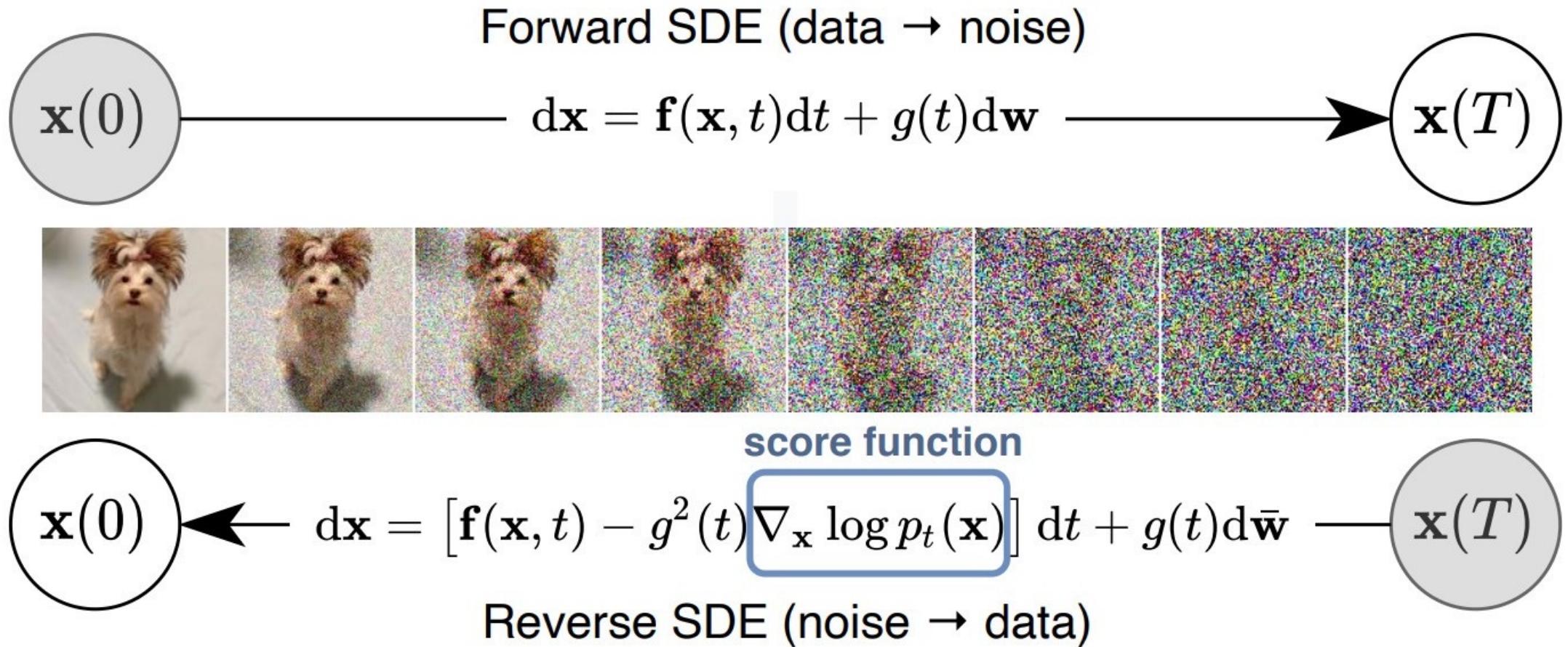
- (+) High quality
- (+) Diversity
- (-) Slow

Diffusion Models

The generation process of a diffusion model is a denoising process.



Score-Based Model / DDPM / DDIM



Diffusion Models

- (+) High quality
- (+) Diversity
- (-) Slow
- (+) (Relatively) easy to implement and train
- (+) Easy to convert a conditional model
- (+) Easy to personalize
- (+) Easy to distill knowledge
- ...

Leveraging Class Labels or Prompts

Classifier Guidance / Classifier-Free Guidance



Conditional Generation

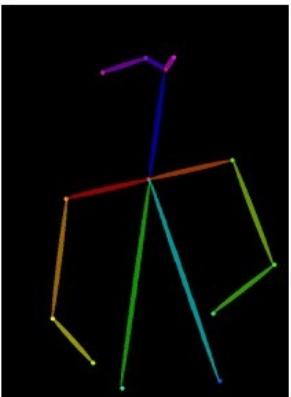
ControlNet



Input Canny edge



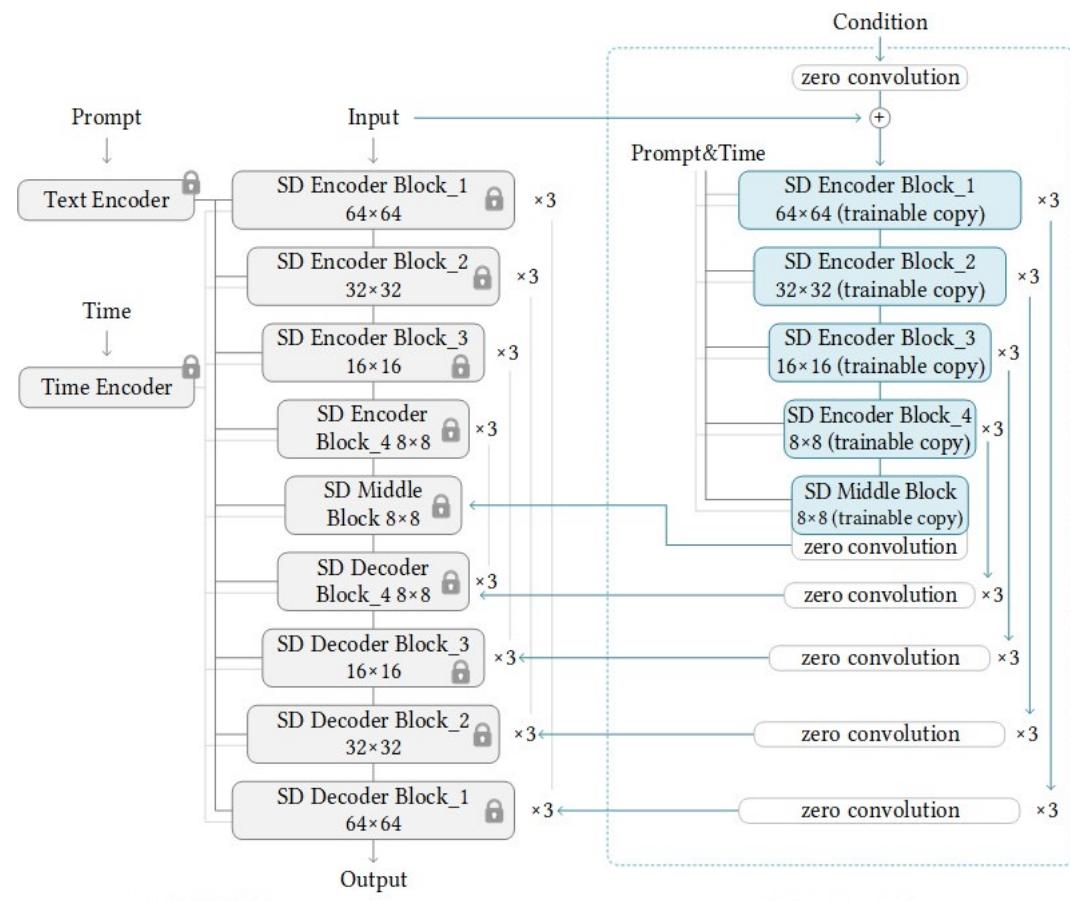
Default



Input human pose



Default

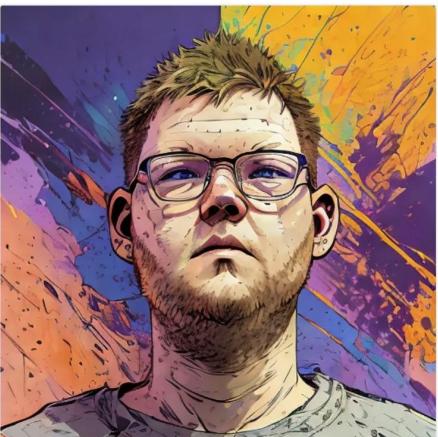


(a) Stable Diffusion

(b) ControlNet

Stylization

LoRA



Personalization

DreamBooth

Input images



A [V] backpack in
the Grand Canyon



A [V] backpack with
the night sky



A [V] backpack in
the city of Versailles



A wet [V] backpack
in water

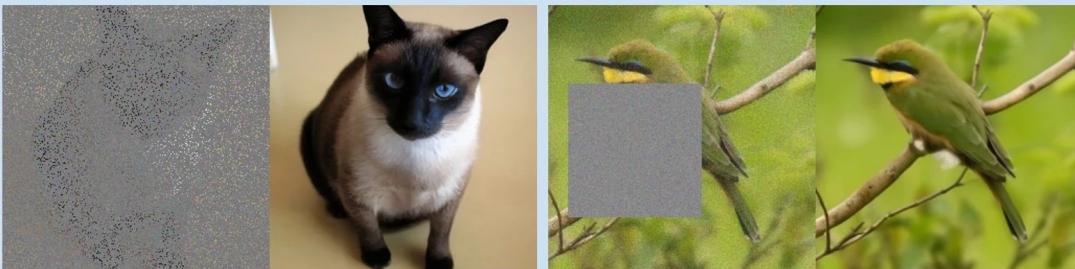


A [V] backpack in Boston

Inverse Problems

Diffusion Posterior Sampling

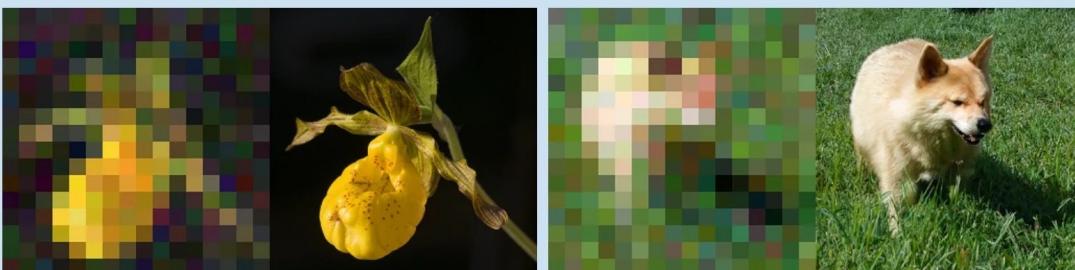
(a) Inpainting



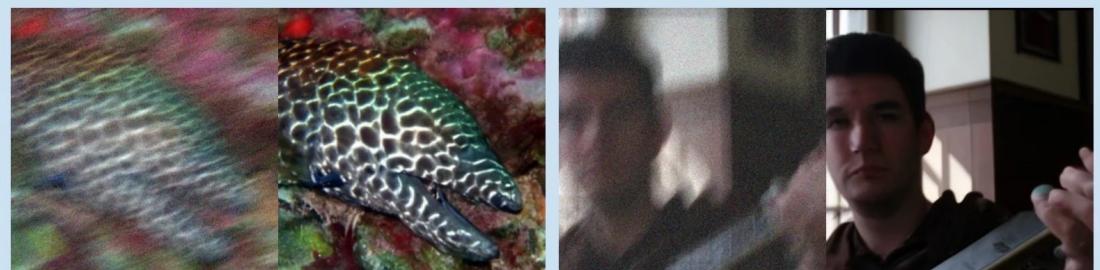
(c) Gaussian deblur



(b) Super-resolution



(d) Motion deblur

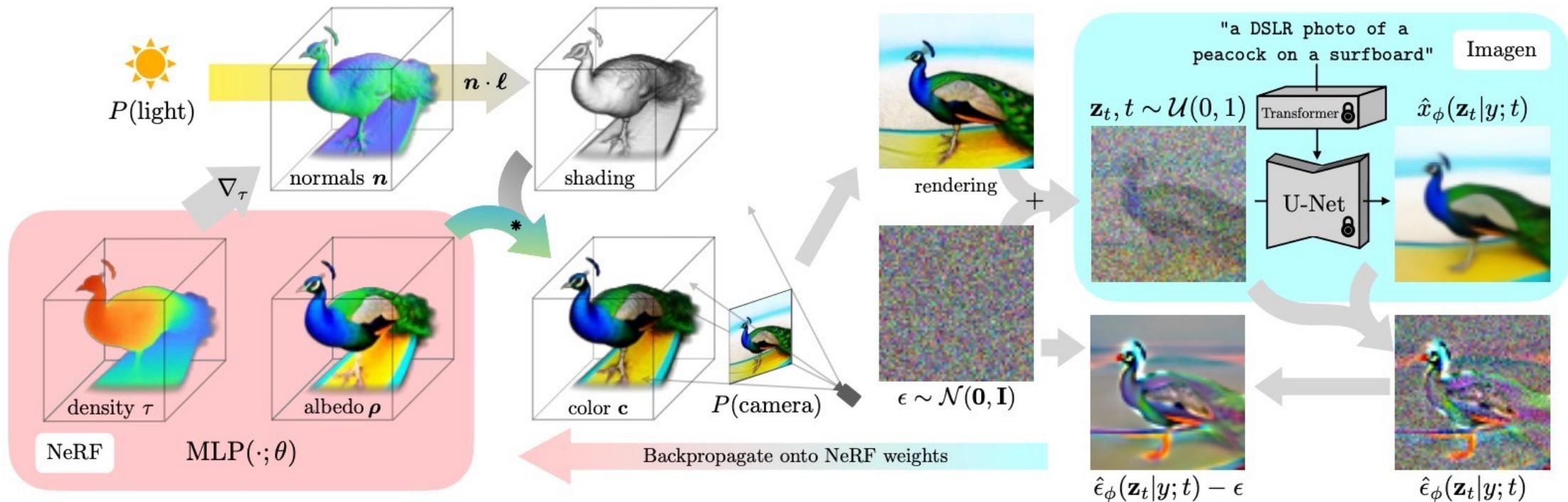


Video Editing



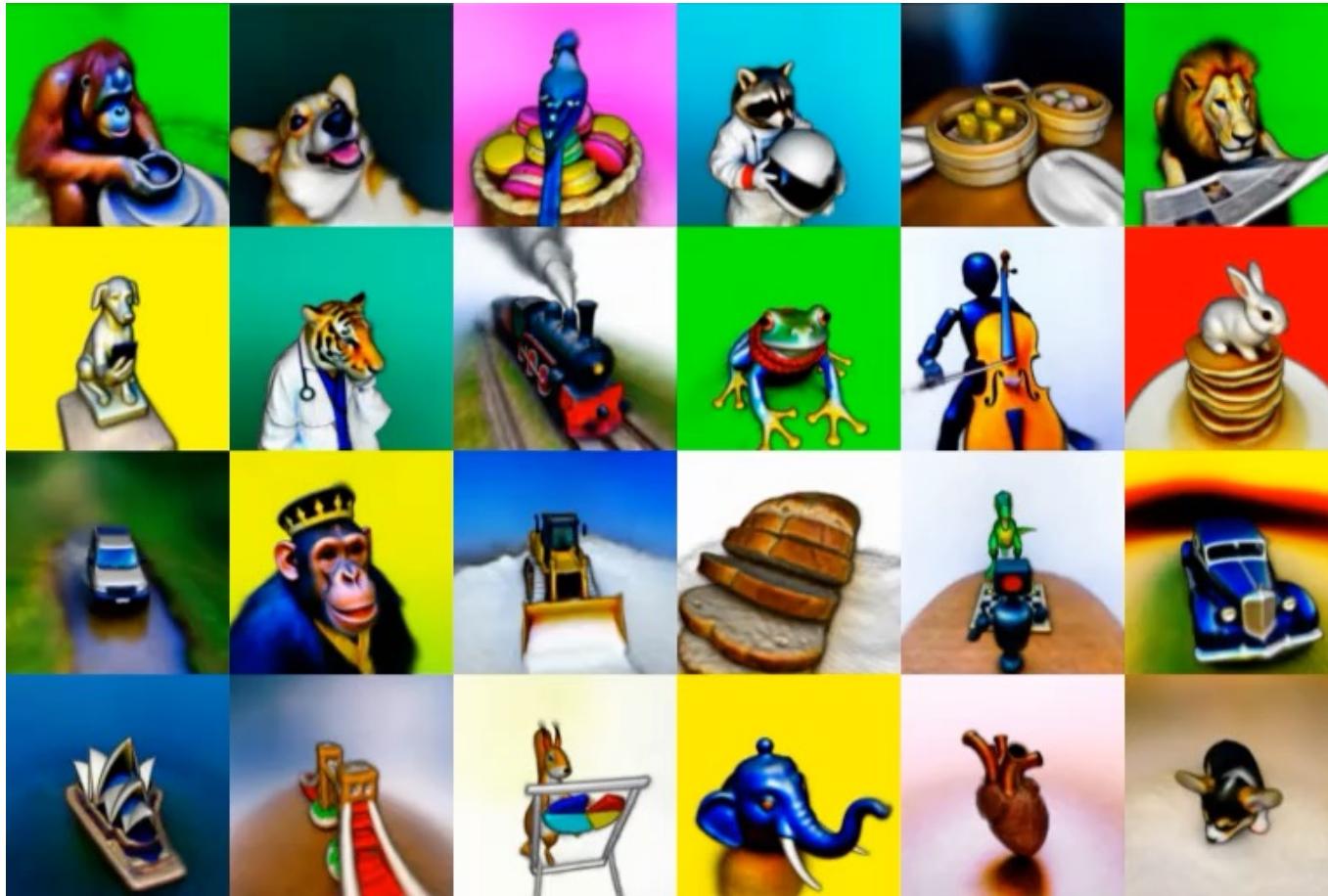
Knowledge Distillation

Score Distillation Sampling



Knowledge Distillation

Score Distillation Sampling



Knowledge Distillation

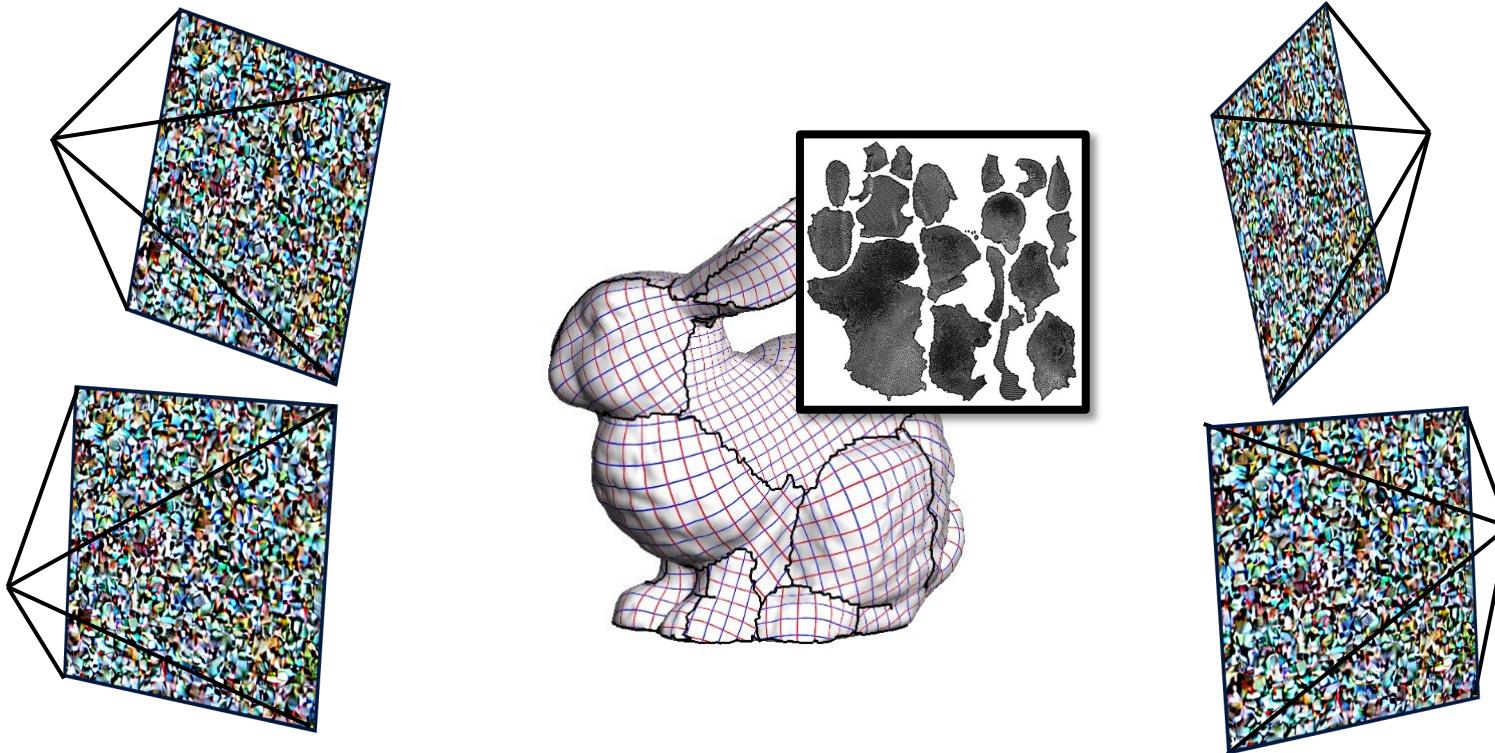
Posterior Distillation Sampling



“Leonardo DiCaprio”



Joint Denoising



Joint Denoising

SyncDiffusion



Joint Denoising

SyncTweedies



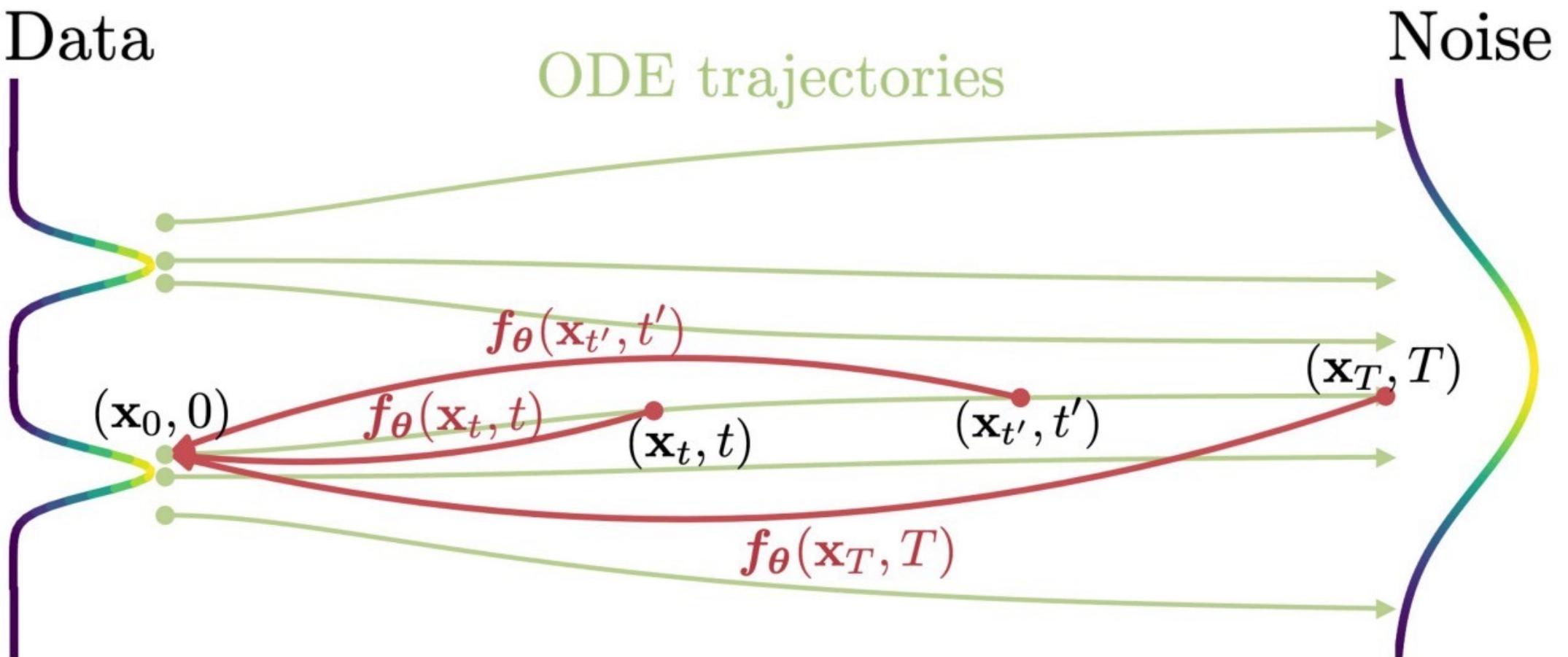
"A Chinese style lantern"



"A car with graffiti"

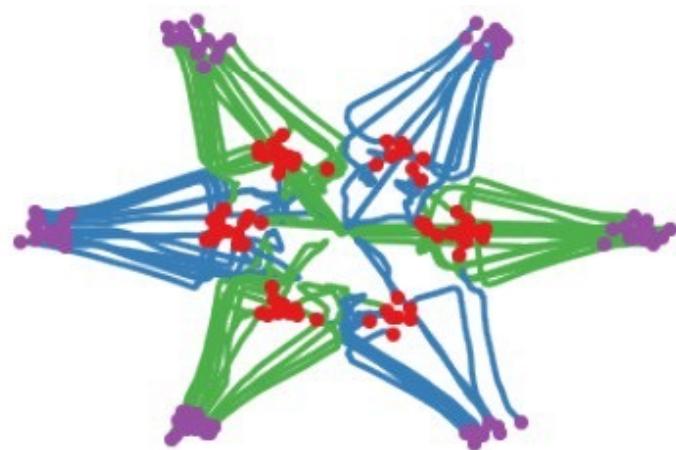
Few-Step Generation

Consistency Models



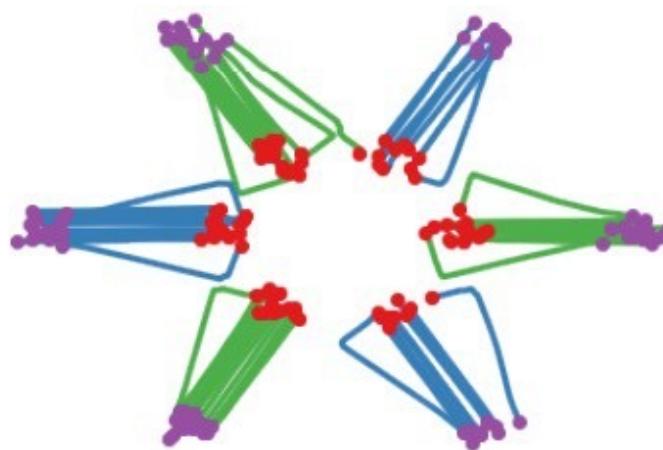
Rectified Flow

ReFlow



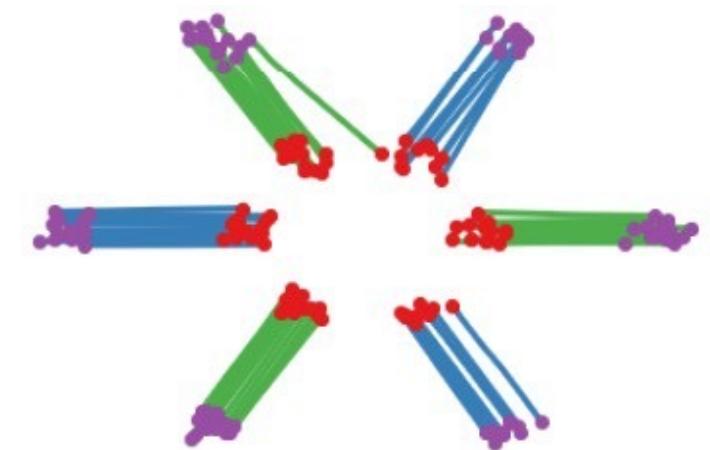
(a) The 1st rectified flow Z^1

$$Z^1 = \text{RectFlow}((X_0, X_1))$$



(b) Reflow Z^2

$$Z^2 = \text{RectFlow}((Z_0^1, Z_1^1))$$



(c) Reflow Z^3

$$Z^3 = \text{RectFlow}((Z_0^2, Z_1^2))$$

Topics

- Background of Generative Models
- DDPM / DDIM / Score-Based Models
- CFG / Latent Diffusion
- Conditional Generation
- Stylization / Personalization
- Inverse Problem
- Knowledge Distillation
- Diffusion Synchronization
- SDE/ODE Solvers
- Consistency Models / Flow-Based Models
- DiT / Applications / Future of Generative Models

Resources

KAIST, Fall 2024, CS492(D): Diffusion Models and Their Applications

<https://mhsung.github.io/kaist-cs492d-fall-2024/>

CS492(D): Diffusion Models and Their Applications

Minhyuk Sung, KAIST, Fall 2024

The collage displays a variety of images generated by a diffusion model, illustrating its capabilities in generating diverse visual content. The images include:

- A wide-angle photograph of a city skyline at dusk or night, with lights reflecting on water in the foreground.
- A highly stylized, colorful illustration of a dog's head, featuring large, expressive ears and a white face with orange and pink accents.
- A painting of a large truck, rendered in a vibrant, abstract style with heavy brushstrokes.
- A photograph of a small potted plant with a single red flower.
- A photograph of a man sitting on a white chair against a plain wall.
- A small, detailed illustration of a sea turtle.
- A green trash bin overflowing with trash bags.
- A classic red armchair with gold-colored legs.
- A potted plant with autumn-colored leaves (orange, yellow, and red).

a painting of a truck

Resources

SIGGRAPH 2024 Course: Diffusion Models for Visual Content Generation

https://geometry.cs.ucl.ac.uk/courses/diffusion4ContentCreation_sigg24/

Diffusion Models for Visual Content Generation

time t

1000 750 500 250 0

Siggraph 2024 Course

July 28 - August 1, Denver, USA

Niloy J. Mitra Duygu Ceylan Or Patashnik Danny CohenOr Paul Guerrero Chun-Hao Huang Minhyuk Sung
UCL/Adobe Adobe Tel-Aviv University Tel-Aviv University Adobe Adobe KAIST

Resources

Tutorials

- CVPR 2023: Denoising Diffusion Models: A Generative Learning Big Bang.
- CVPR 2024: 3D/4D Generation and Modeling with Generative Priors.
- CVPR 2024: Diffusion-based Video Generative Models.
- CVPR 2025: From Video Generation to World Model.

Blogs

- "Generative Modeling by Estimating Gradients of the Data Distribution", Yang Song.
- "What are Diffusion Models?", Lilian Weng.
- "Understanding Diffusion Models: A Unified Perspective". Calvin Luo.
- "Tutorial on Diffusion Models for Imaging and Vision". Stanley H. Chan.
- "Step-by-Step Diffusion: An Elementary Tutorial". Preetum Nakkiran, Arwen Bradley, Hattie Zhou, and Madhu Advani.

In this course...

We will discuss diffusion models, covering both their
theoretical foundations and practical applications.

Prerequisite

- **Background in machine learning/deep learning.** We'll specifically focus on diffusion models (while briefly discussing the background of generative models).
- **Experience with neural network implementation.** There will be programming assignments and a final project.
- Recommended prior courses:
 - 線性代數
 - 機率
 - 微分方程
 - 深度學習+深度學習實驗

Week	Date	Topic	Assignments
1	2025-09-04	GAN / VAE	
2	2025-09-11	DDPM	#1 – DDPM
3	2025-09-18	DDIM	
4	2025-09-25	CFG / Latent Diffusion / ControlNet / LoRA / Zero-Shot Applications	#2 – DDIM & LoRA
5	2025-10-02	DDIM Inversion / Score Distillation (attending Corl 2025)	
6	2025-10-09	Diffusion Synchronization / Inverse Problems	#3 – Distillation
7	2025-10-16	Probability Flow ODE / DPM-Solver	
8	2025-10-23	Flow Matching (attending ICCV 2025)	#4 – Flow Matching
9	2025-10-30		
10	2025-11-06	Final Project Proposal	
11	2025-11-13		
12	2025-11-20		
13	2025-11-27	Paper Presentation	
14	2025-12-04	Guest Lecture – Ta-Ying (Tim) Cheng	
15	2025-12-11	Guest Lecture – Chieh (Hubert) Lin	
16	2025-12-18	Guest Lecture – Chih-Hao Lin	
17	2025-12-25	Final Project Presentation	

Lectures

- 100 mins IN-PERSON lecture + 50 mins ONLINE recorded video.
- Recordings will be available on E3.

Guest Lectures



Ta-Ying Cheng
Research scientist
@ Netflix
Topic: Material editing



Chieh Hubert Lin
Research scientist
@ Architectural design startup
Topic: 3D inpainting



Chih-Hao Lin
3rd year PhD student
@ UIUC
Topic: Weather Synthesis

Course Logistics

Yu-Lun Liu (Instructor)

- Assistant Professor in the Department of Computer Science.
- Homepage: <https://yulunalexliu.github.io/>
- Office: EC713.

Teaching Assistants



黃怡川 (資科工所博三)
yichuanh.cs12@nycu.edu.tw



李杰穎 (資科工所博一)
jayinnn.cs14@nycu.edu.tw



陳映寰 (資科工所碩二)
yinghuachen.cs13@nycu.edu.tw



張欗齡 (資科工所碩二)
siang1105.cs13@nycu.edu.tw



鄭淮薰 (資科工所碩二)
huaish.cs13@nycu.edu.tw



柯柏旭 (資科工所碩二)
hentci.cs13@nycu.edu.tw

Evaluation

- Programming Assignment: 40% (10% each)
- Paper Presentation: 15%
 - 10%: Hacker's Deliverables (code, demo, presentation)
 - 5%: 問問題
- Final Project: 40%
 - 10% Proposal Presentation
 - 15% Oral & Demo Presentation
 - 10% Final Report (4 pages in CVPR LaTeX format)
 - 5%: 問問題
- Participation in Guest Lectures: 5%

Programming Assignments

- Assignment1 (DDPM)
- Assignment2 (DDIM & LoRA)
- Assignment3 (Distillation)
- Assignment4 (Flow Matching)

Prerequisite

- You'll need **basic programming skills in Python and PyTorch** to complete the programming assignments.
- **Start the programming assignments as early as possible!**

Submission

- Each programming assignment is due **two weeks after** the assignment session.
- Submit your solutions on E3.
- **No late submissions allowed.** I.e., you will get zero credit even with only 1 minute late submitting deliverables on E3.

Paper Presentation

- 每組四位學生（請現在就開始找組員！）
- 報告的組別：報告近三年發表於
CVPR/ICCV/ECCV/NeurIPS/ICLR/ICML/SIGGRAPH/SIGGRAPH Asia main conference 的論文
 - 根據你報告的好壞我會給予評分（總成績 5%）
 - **Hacker's Deliverables**（總成績 5%）
- 其他所有組別：問問題
 - 在 presentation 中間或結束後皆可發問
 - 每個 presentation 我會選出最 insightful 的問題加分（每次總成績 1%）

Paper Presentation – Hacker

報告的組別必須扮演一個需要盡快展示這篇論文的駭客

- 在你自己的資料/任務上展示這篇論文的結果
- 準備與班上同學分享演算法的核心程式碼並展示你的實作
- 不要只是下載並執行現有的實作



Paper Presentation – Hacker



把公開的 Github repo
的 README 跑完，展
示 paper 上就有的結果



- 用自己的資料/prompt 去跑
- 稍微改一點點 code 讓 visualization 更好看
- 搭配論文內容講解對應的 程式碼區塊
- 展示 code diff 或說明安裝 環境踩到的坑

Paper Presentation – 其他組別問題



不著邊際非常飄渺的問題：“這篇 paper 的東西能用在 video 上嗎？”



- “為什麼不用 X ，要這麼複雜？好處是 Y 嗎？”
- “用 Z 演算法以及 A 改動，就可以 training-free 了嗎？”
- “這樣做會有 B 的問題，也許可以使用 C 方法解決，因為...”

Final Project

- 每組四位學生（與 Paper Presentation 組別相同）
- 三種選項：
 - [量變] 改善現有論文的演算法（生成品質、速度等等）
 - [質變] 創意開發新應用或演算法（影像 -> 視訊, training -> zero-shot）
 - [復現] Re-implement 近五年任何一篇沒有公開程式碼的頂尖會議論文

Final Project Oral Presentation

- 報告的組別：
 - 根據創意/實用度/完整度/報告的好壞等等我會給予評分（總成績 15%）
- 其他所有組別：扮演 reviewer 問問題
 - 在 presentation 中間或結束後皆可發問
 - 每個 presentation 我會選出最 insightful 的問題加分（每次總成績 1%）



Final Project Final Report

四頁的 paper in CVPR LaTeX format

- Title
- Abstract
- Introduction
- Related Work
- Method
- Experiments
- Conclusion
- References



AI Coding Assistant Tool Policy

- 你可以在 programming assignments 和 final project 中使用 AI coding tools，例如 ChatGPT、Copilot、Codex 和 Code Intelligence
- 但直接抄網路上或別人的程式碼仍然是嚴格禁止的
 - 學期成績零分，並向學校報告

Plagiarism

- 我們將使用 Turnitin 與 Moss 來比對原創性檢測與遏止抄襲
- Turnitin: <https://www.turnitin.tw/>
- Moss: <https://theory.stanford.edu/~aiken/moss/>

Computing Resource

本課程**不提供**任何運算資源，請使用個人/實驗室/付費雲端資源完成 programming assignments 與 final project

你會需要至少一張有 12 GB VRAM 的 NVIDIA GPU