# A Quick Tour of Computer Vision
## SAST Summer Training 2023
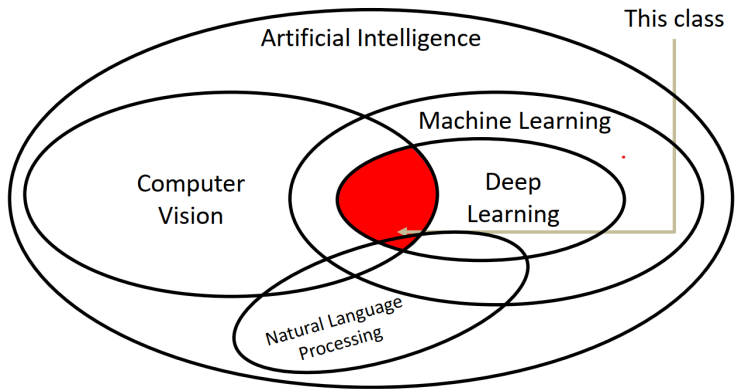
Kevin Zhang

THU CST

July 30th, 2023

## Today's Goal

- Explore the applications of previously learned neural networks (e.g. CNN, Transformers) in basic computer vision tasks.
- Understand the fundamentals of cutting-edge generative models: VAE, GANs, and Diffusion Models.
- Explore the widely used open-source diffusion model codebase: *StableDiffusion*.
- Develop your interest and spark your inspiration :)

## What's Computer Vision?

**Wiki**: Computer vision tasks include methods for **acquiring, processing, analyzing and understanding** digital images, and extraction of high-dimensional data from the real world in order to produce numerical or symbolic information.

**Note**: Computer Vision $\not\subset$ Deep Learning

# What's Computer Vision?

# What's Computer Vision?



Face detection / recognition



Self-driving cars



Human care



Medical image analysis



Remote sensing / earth observing

## What's Computer Vision?

A WIDE range of computer vision tasks:

- Different modalities
  - Image
  - Video
  - 3D Object/Scene
  - Multi-modal (Vision, Language, Speech...)
- Different targets
  - Recognition
  - Segmentation
  - Detection
  - Stylization
  - Captioning
  - Generation
  - ...

## Today's Focus

A WIDE range of computer vision tasks:

- Different modalities
    - Image → *Basic*
    - Video
    - 3D Object/Scene
    - Multi-modal (Vision, Language, Speech...) → *Exciting*

- Different targets
    - Recognition → *Basic*
    - Segmentation
    - Detection
    - Stylization
    - Captioning
    - Generation → *Exciting*
    - ...

**1** Introduction

**2** Basic Task: Image Classification
  Introduction
  CNN for Image Classification
  Vision Transformer (ViT): Towards Larger Model

**3** Recent Progress: Generation Models

**4** Reference

**1** Introduction

**2** Basic Task: Image Classification
Introduction
CNN for Image Classification
Vision Transformer (ViT): Towards Larger Model

**3** Recent Progress: Generation Models

**4** Reference

## Why Image Classification?

- Core computer vision task: make computer perceive the world
- Building blocks for other tasks: detection, captioning, and even AlphaGo...
- The first challenge that deep learning made great success over classical methods
- A good chance to review and get a deeper understanding of previously learned neural network architectures, such as CNN and Transformers

## What's Image Classification?



This image by Nikita is licensed under CC-BY 2.0

(assume given a set of labels)
{dog, cat, truck, plane, ...}

$\longrightarrow$

**cat**
dog
bird
deer
truck

## Challenge for Image Classification

**The Problem**: Semantic Gap



This image by Nikita is licensed under CC BY 2.0
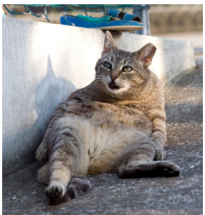


What the computer sees

An image is a tensor of integers between [0, 255]:

e.g. 800 x 600 x 3
(3 channels RGB)

## Challenge for Image Classification

So many different cats!

## Challenge for Image Classification

So many different cats (or cat tails)!

# Opportunity: Big data



IMⱯGENET **Large Scale Visual Recognition Challenge**

The Image Classification Challenge:
1,000 object classes
1,431,167 images

Output:
Scale
T-shirt
Steel drum
Drumstick
Mud turtle

Deng et al, 2009
Russakovsky et al. IJCV 2015

# Imagenet ILSVRC



ImageNet Classification top-5 error (%)

Introduction
○○○○○○○

Basic Task: Image Classification
○○○○○○○○○○●○○○○○○○○○○○

Recent Progress: Generation Models
○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○

Reference
○○○

## Recap: CNN architecture



Scan
the image

Generate
hierarchy of features

Recognize from
high level features

Extract features: convolution layers + pooling layers + normalization layers

## Recap: Local Connectivity in CNN



Key assumption 1 (Locality): Assume local information is enough for feature extraction and recognition.

# Recap: Parameter Sharing in CNN



patch    pixel level inputs    $1^{st}$ level features    $2^{nd}$ level features

$\begin{array}{|c|c|}\hline \theta_{11} & \theta_{12} \\\hline \theta_{13} & \theta_{14} \\\hline\end{array}$ weights

Receptive field

$\begin{array}{|c|c|}\hline \theta_{21} & \theta_{22} \\\hline \theta_{23} & \theta_{24} \\\hline\end{array}$ weights

Key assumption 2 (Shift invariance): If a feature is useful at one position, then it should also be useful at other positions.

# AlexNet (2012)



first use of ReLU

heavy data augmentation

7 CNN ensemble: 18.2% → 15.4%

Parameter Inefficient!

# GoogLeNet (Inception, 2014): Going deeper

Deeper networks, with computational efficiency



**Stem Network**

**Stacked Inception Modules**

**Classifier**

- 22 layers
- Efficient Inception modules
- The first one without FC layers
- Only 5 million parameters! 12x less than AlexNet
- ILSVRC'14 classification winner (6.7% top 5 error)

# ResNet (CVPR best paper, 2015): Going really deep

Residual Block



Plain layers          Residual block

- Key observation: Deeper models are hard to optimize.
- Idea:
  - Copying the learned layer from the shallow model
  - Setting additional layers to identical mapping
- Solution: Instead of learning $H(x)$ directly, learn the residual $F(x) = H(x) - x$.

# ResNet (CVPR best paper, 2015): Going really deep

- ResNet architecture:

  - Stack residual blocks

  - Every residual block has two 3x3 conv layers

  - Periodically, double #filters and downsample spatially using stride 2 (/2 in each dimension)

  - Additional conv layer at the beginning

  - Global average pooling at the end (FC layer only to output classes)



Residual Block

# Well-known CNN models that emerged between 2012 and 2017



Top1 vs. Architectures



Top1 vs. Operations, Size ∝ Parameters

搭积木？知乎：像 ResNet、SENet 这些网络是怎么想出来的？

**1** Introduction

**2** Basic Task: Image Classification

Introduction

CNN for Image Classification

Vision Transformer (ViT): Towards Larger Model

**3** Recent Progress: Generation Models

**4** Reference

# Recap: Transformer



- Main Components:
  - Scaled Dot-product Attention
  - (Masked) Multi-head Attention
  - Position-wise FFN
  - Residual Connections
  - Layer Normalization
  - Positional Encoding
- Architecture:
  - Encoder →
    *A new way to extract features*!
  - Decoder with Masking
  - Encoder-Decoder Attention

# ViT: An image is worth 16*16 words

Simple idea: Split the image into fixed-size patches and treat the image as a sequence of patches! (The same Transformer Encoder architecture as before)

$$\mathbf{x}^{H \times W \times C} \rightarrow \mathbf{x}_p^{N \times (P^2 \cdot C)} \rightarrow \mathbf{z}^{N \times D}$$



**Vision Transformer (ViT)**

**Transformer Encoder**

## ViT: An image is worth 16*16 words

- Less inductive bias $\rightarrow$ One architecture for **multi-modal data** and **multiple downstream tasks**!

- Easy to **scale up**! (Demanding for extremely large scale dataset, model regularization and data augmentation)

**1** Introduction

**2** Basic Task: Image Classification

**3** Recent Progress: Generation Models

   Transformer-based Auto-regressive Models
   Generative Adversarial Networks (GANs)
   Diffusion Models

**4** Reference

Introduction
○○○○○○○○

Basic Task: Image Classification
○○○○○○○○○○○○○○○○○○○○○○○○○

Recent Progress: Generation Models
○●○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○

Reference
○○○

# AIGC: A new opportunity

**1** Introduction

**2** Basic Task: Image Classification

**3** Recent Progress: Generation Models
Transformer-based Auto-regressive Models
Generative Adversarial Networks (GANs)
Diffusion Models

**4** Reference

## Transformer again: Unify texts and images

Challenges for applying GPT-like large models to image generation:

- Huge computation consumption if generating in pixel space (A single image might have tens of thousands of pixels)
- Two-dimensional information is probably ignored in sequence generation

Initial solution: Reduce dimensions and generate in the latent space

## Generative Modeling

Learn the probability distribution $p_\theta(x)$ that generates the data.



Data: $\{(x_i)\}_{i=1}^n$ → Generative Modeling $f(x)$ → Distribution: $P_\theta(x)$

# Variational Auto-Encoder (VAE)



Sample $x$ from $x|z \sim N(\mu_{x|z}, \Sigma_{x|z})$

$\mu_{x|z}$    $\Sigma_{x|z}$

Decoder network
$p_\theta(x|z)$

$z$

Sample $z$ from $z|x \sim N(\mu_{z|x}, \Sigma_{z|x})$

$\mu_{z|x}$    $\Sigma_{z|x}$

Encoder network
$q_\phi(z|x)$

$x$

Training:

- **Encoder** output should match the prior distribution, e.g. $\mathcal{N}(0, 1)$.
- **Decoder** output should reconstruct the input distribution.

Inference:

- Throw the encoder away, sample from the prior distribution and use the decoder to get output image.

# Variational Auto-Encoder (VAE)

## You want more math? (For those interested, not required)

Goal: Maximize $p_\theta(x) \rightarrow$ Objective: Evidence lower Bound, ELBO
($\theta$: Decoder parameters, $\phi$: Encoder parameters)

$$\log p_\theta(x_i) = \mathbb{E}_{z \sim q_\phi(z|x_i)}[\log p_\theta(x_i)]$$
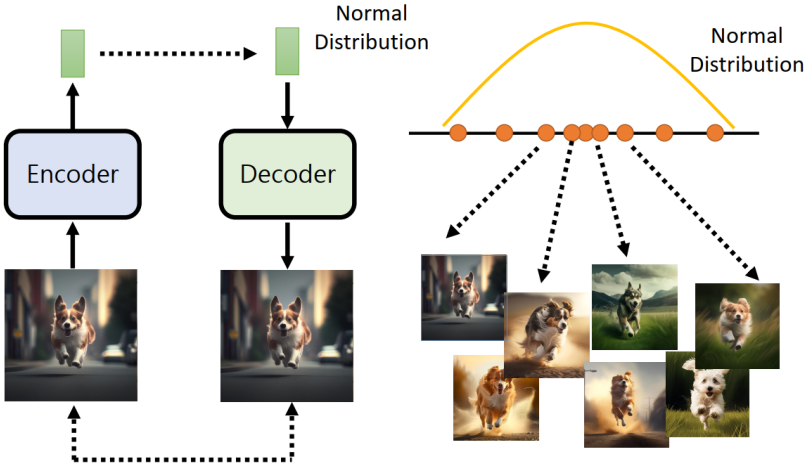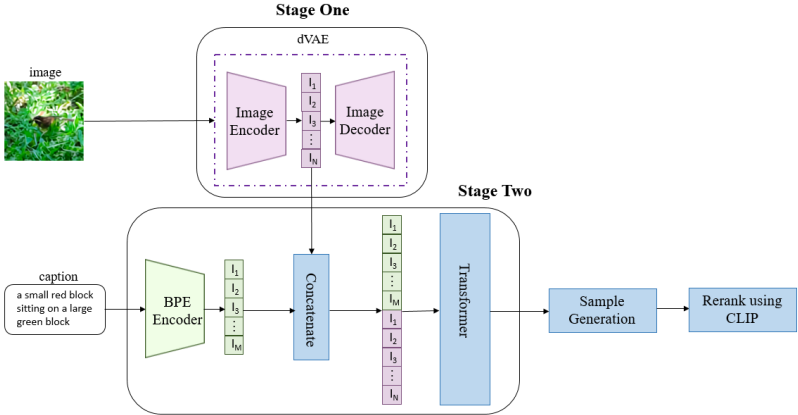
$$= \mathbb{E}_z[\log \frac{p_\theta(x_i|z)p_\theta(z)}{p_\theta(z|x_i)}]$$

$$= \mathbb{E}_z[\log p_\theta(x_i|z)\frac{p_\theta(z)}{q_\phi(z|x_i)}\frac{q_\phi(z|x_i)}{p_\theta(z|x_i)}]$$

$$= \mathbb{E}_z[\log p_\theta(x_i|z)] - \mathbb{E}_z[\log \frac{q_\phi(z|x_i)}{p_\theta(z)}] + \mathbb{E}_z[\log \frac{q_\phi(z|x_i)}{p_\theta(z|x_i)}]$$

$$= \mathbb{E}_z[\log p_\theta(x_i|z)] - KL(q_\phi(z|x_i)||p_\theta(z))$$

$$+ KL(q_\phi(z|x_i)||p_\theta(z|x_i))$$

$$\geq \mathbb{E}_z[\log p_\theta(x_i|z)] - KL(q_\phi(z|x_i)||p_\theta(z))$$

# DALL-E (OpenAI)

Introduction
○○○○○○○

Basic Task: Image Classification
○○○○○○○○○○○○○○○○○○○○○○○

Recent Progress: Generation Models
○○○○○○○○○○●○○○○○○○○○○○○○○○○○○○○○○

Reference
○○○

# CogView (Tsinghua)



The framework of CogView. [ROI1], [BASE1], etc., are seperator tokens.

## Pros and cons of VAE

Pros:

- Principled approach to generative models with solid mathematics basis.
- Allow inference of encoder $q_\phi(z|x_i)$, can be useful low-dimension feature representation for other tasks and models.

Cons:

- Evidence Lower bound: an over-simplified approximation.
- Intractable for complex distributions.
- Blurrier and lower quality than GANs and Diffusion models.

# GAN: A two-player game

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

Discriminator wants
D(x) = 1 for real data

Discriminator wants
D(x) = 0 for fake data

$$\min_{G} \max_{D} \left( E_{x \sim p_{data}}[\log D(x)] + E_{z \sim p(z)} \left[ \log \left( 1 - D(G(z)) \right) \right] \right)$$

Generator wants
D(x) = 1 for fake data



Sample z from $p_z$ → z → Generator Network G → Generated Sample → Discriminator Network D → Fake / Real

- The center of GAN is an adversarial loss! $\rightarrow$ An idea easy to generalize.

## Conditional GAN



One-hot label vector $y$

Control

Each row shares same label

- Provide condition information to indicate what to generate.

# StyleGAN



- Latent $\mathbf{z} \to$ style features, noise $\to$ content features (diverse details).

# DragGAN (SIGGRAPH 2023)



- Utilizing the continuous and discriminative latent space of *StyleGAN* to apply natural and meaningful editing.
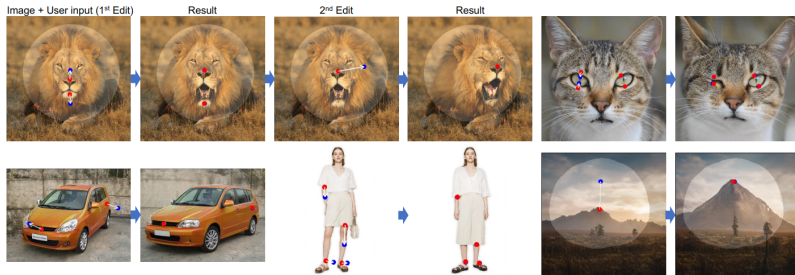
## Pros and cons of GAN

Pros:

- High-quality generated images (enough to fool the discriminator and also your eyes!)
- High sampling efficiency!
- Relatively continuous latent space (which stores meaningful semantic information).

Cons:

- Hard to train! (A minimax game) $\rightarrow$ Relatively difficult to scale up.
- Mode collapse (Relatively low diversity) $\rightarrow$ A lot of follow-up works to mitigate the problem, e.g. *WGAN*.
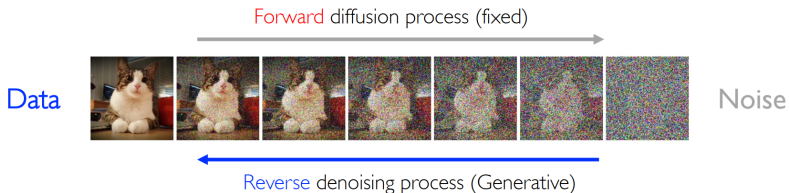
**1** Introduction

**2** Basic Task: Image Classification

**3** Recent Progress: Generation Models
    Transformer-based Auto-regressive Models
    Generative Adversarial Networks (GANs)
    Diffusion Models

**4** Reference

## Denoising Diffusion Probabilistic Model (DDPM)



Forward diffusion process (fixed)

Data

Noise

Reverse denoising process (Generative)

Denoising diffusion models consist of two processes:

- **Forward** **diffusion process**: Adds noise to input gradually.
- **Reverse** **denoising process**: learns to generate by denoising.

## Training



$x_0$: clean image



$\varepsilon$: noise

**Algorithm 1** Training

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

Noisy image

$\bar{\alpha}_1, \bar{\alpha}_2, \ldots \bar{\alpha}_T$
─────────────
smaller

Target
Noise

Noise
predictor

## Training



$x_0$

$\varepsilon$

Sample $t$

$\sqrt{\bar{\alpha}_t}$   $x_0$   $+ \sqrt{1 - \bar{\alpha}_t}$   $\varepsilon$   $=$

Noise Predicter

t

?????

$\varepsilon$

# Training



What you imagine...

Random sample

Step 1    Step 2    input

input

ground truth

Actually...

$$\sqrt{\bar{\alpha}_t} \quad + \sqrt{1 - \bar{\alpha}_t} \quad = $$

$x_0$    $\varepsilon$    input

ground truth

# Inference

## You want more math? (For those interested, not required at all)

- Lil's blog: What are Diffusion Models?
- 知乎：扩散模型之 DDPM
- My own notes (uploaded in Tsinghua Cloud)

## VAE v.s. DDPM



**VAE**    Encoder → ▮ → Decoder

**Diffusion**    Add noise / X N →   → Denoise   (X N)

- One step v.s. $N$ steps: Efficiency & Quality tradeoff
- Latent low-dimension representation v.s. Gaussian noise of the same shape

## Stable Diffusion: Conditional diffusion in latent space



- Try it on huggingface: https://huggingface.co/spaces/stabilityai/stable-diffusion

Something practical: Stable Diffusion open-source codebase

- Github repo:
  https://github.com/CompVis/stable-diffusion
- Built on PyTorch Lightning
- Widely used by mountains of downstream tasks

## Usage: Inference

1. Create conda environment:

```
1  conda env create -f environment.yaml
2  conda activate ldm
```

2. Download checkpoints from huggingface and put or link it in /models folder:

```
1  mkdir -p models/ldm/stable-diffusion-v1/
2  ln -s <path/to/model.ckpt>
3      models/ldm/stable-diffusion-v1/model.ckpt
```

3. Edit config file in the /config folder

4. Sample images:

```
1  python scripts/txt2img.py --prompt "..."
```

## Usage: Training

```
1  python main.py -t -b <config-files> -l <log-file>
```

(The preceding procedure is the same as inference)

## Code structrue

**1** Introduction

**2** Basic Task: Image Classification

**3** Recent Progress: Generation Models

**4** Reference

Reference and further reading

- All the reference papers are embedded as a hyperlink in the corresponding slide title.
- Some of the pictures in this slide are credited to the following wonderful courses:
    - Stanford University CS231n: Deep Learning for Computer Vision
    - National Taiwan University: Machine Learning by Hung-Yi Lee
    - Tsinghua University: Deep Learning by Mingsheng Long

*Thanks!*