# THU×SENSETIME – 80231202

# Advanced Computer Vision

Friday, February 26, 2021

# Content

# Course Introduction

## Overview

This course involves **computer vision**, **deep learning** and other fields of knowledge. It elaborates with the latest academic achievements and practical cases of industrial scenes and explain the classic and state-of-the-art methods in computer vision.

| What we have | What you will learn | What you need |
|---|---|---|
| • Focus on Both Classics and Frontiers<br><br>• Combination of Academia and Industry<br><br>• Teaching from the shallower to the deeper<br><br>• GPU clusters for experiments | • Basic theories and advanced methods in Computer Vision<br><br>• Understand and explore practical problems in the industry<br><br>• Improve your research ability and innovative ability | • **Mathematics**<br> • Calculus<br> • Linear Algebra<br> • Basic Probability and Statistics<br>• **Coding ability**<br> • **Python** is recommended<br>• Machine Learning |

# Course Introduction

## Chapter 1

- Basics of computer vision & image processing
- Introduction of the neural network and deep learning framework

## Chapter 2

- cutting-edge research directions in computer vision
- the algorithm model optimization and performance improvement methods in visual scenes.

## Chapter 3

- the practical problems faced by computer vision and the solution ideas in combination with the specific scenes of industry.

# Syllabus

## Chapter 1 - Computer Vision Overview and Deep Learing Basics

1. Computer Vision Basics
2. Image and Video Processing
3. Feature Detection
4. CNN & High-level Feature Extraction
5. Training Framework and Model Optimization

## Chapter 2 - Advanced Computer Vision Tasks

6. Image Classification
7. Object Detection
8. Image Segmentation
9. Video Understanding and Sequence Analysis
10. Low-Level Computer Vision Task
11. Neural network Model Acceleration and Compilation
12. 3D Vision
13. Representation Learning in Vision Tasks

## Chapter 3 - Lectures on industry applications

14. Smart City
15. AutoPilot
16. 3D Vision and Augmented Reality

# Course Introduction

- **Textbook**
  - ***Computer Vision Algorithms and Applications***
    - by Richard Szeliski
    - Preview version：[Link]
  - ***Pattern Recognition and Machine Learning***
    - by Christopher Bishop
    - Free online version：[Link]
  - ***Deep Learning***
    - by Goodfellow, Bengio, and Courville
    - Index：[Link]
  - ***Dive into deep learning***
    - An interactive deep learning book with code, math, and discussions, based on the NumPy interface
    - Free online version：[Link]

# Course Introduction

- **Grading Policy**

**Quizzes (20%)**

2 Quizzes in class, completed by one person

**Assignments (30%)**

3 Assignments finish after class by one person

**Final Project (50%)**

Choose one topic , submit 3-4 pages paper and make a oral presentation during the seminar.

Collaboration in groups of up to three people.

- **Assignment**
  - All assignments should be finished by one person
  - You can finish assignment on your local machines or on clusters provided by SenseTime
  - More details will be update on Course Homepage

| Assignment | Released Date | Due Date | Topic |
|---|---|---|---|
| Assignment 1 | Mar. 12 | Mar. 26 | Deep learning training framework and model optimization implementation |
| Assignment 2 | Apr. 2 | Apr. 16 | Advanced Computer Vision Tasks |
| Assignment 3 | May. 7 | May. 14 | Lightweight Model Quantization and Model Pruning |

# Course Introduction

- **Final Project**
  - Choose one topic and finish the project
  - You should submit
    - One page proposal and discuss it with TAs (topic, idea, method, experiments)
    - A term paper of 4 pages (excluding figures) in maximum, double column, font size is equal or larger than 10
    - Code and sample data
    - Project presentation
  - Collaboration in groups of up to three people

**Mar. 19**

Final Project release

**Apr. 20**

Submit topics of the final project

**Apr. 22**

Tutorial (optional to attend): Discuss with TA-in-charge

**May 9**

Submit proposal (1-2 pages)

**May 28**

Final project Seminar (10 minutes presentation and 3 minutes Q&A)

# Course Introduction

- **Instructors**
  - Dr. Li Yali          liyali13@mail.tsinghua.edu.cn
  - Dr. Dai Jifeng       daijifeng@sensetime.com
  - Dr. Liu Yu           liuyu@sensetime.com
  - Dr. Li Hongyang      lihongyang@sensetime.com
- **TAs**
  - Wang Han            i@hann.wang
  - Wang Cheng          wangcheng@senseauto.com
  - Song Guanglu        songguanglu@sensetime.com
  - Niu Yazhe           niuyazhe@sensetime.com

# Course Introduction

- **Lecture Time & Venue**
  - **Friday**, 9:50am-11:25am
  - **4203**, No.4 Teaching Building

- **Optional Tutorials & QA Time**
  - **Thursday**, 19:00-20:00
  - Tencent Meeting Room：785 271 5223

- **Course Homepage**
  - https://thu-acv.github.io

- **Discussions**
  - WeChat Group
  - Tencent Meeting Room：785 271 5223

THU 高等计算机视觉 课程群

商汤泰坦小助手
中国

**Content**

# What's Computer Vision



Artificial Intelligence (AI)

Machine Learning (ML)

Depp Learning (DL)

Convolutional Neural Network (CNN)

**Computer Vision**

- Object detection
- Object classification
- Scene understanding
- Semantic scene segmentation
- 3D reconstruction
- Object tracking
- Human pose estimation
- Activity recognition
- VQA
- ....

Vision is the most important source of information for the human brain and is the "**entrance hall**" of AI.

**Content**

- **Biological Vision**

- **Ancient Human Vision**

## Camera Obscura

Gemma Frisius, 1545

This work is in the public domain

Encyclopedia, 18th Century

This work is in the public domain

Leonardo da Vinci, 16th Century AD

This work is in the public domain

- **Neuroscience and Vision**



Kanwisher et al. J. Neuro. 1997

Epstein & Kanwisher, Nature, 1998

- **Marr Computational Vision**



3D Reconstruction
Not talent, but
computation

- **Marr Computational Vision**



Stages of Visual Representation, David Marr, 1970s

- **Feature Detection——SIFT**

- **Feature Detection——HOG**



(a)  (b)  (c)  (d)  (e)  (f)  (g)

https://web.archive.org/web/20110408220331/
http://www.acemedia.org/aceMedia/files/document/wp7/2005/cvpr05-inria.pdf

- **3D reconstruction**



Agarwal et al.
ICCV, 2009

- **Image Classification**

- **IMAGENET Challenge**



22,000 categories : 15,000,000 images

J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li & L. Fei-Fei. *CVPR*, 2009.

- **IMAGENET Challenge**

# History of Computer Vision

- ## IMAGENET Challenge

# Image Classification on ImageNet

- **Object Detection**



https://cocodataset.org/

# Object Detection on COCO test-dev

- **Instance Segmentation**



https://www.lvisdataset.org/explore

- **Semantic Segmentation and Instance Segmentation**



Input Image

Semantic Segmentation

Instance Segmentation

# Instance Segmentation on COCO test-dev

# CLIP: Connecting Text and Images

We're introducing a neural network called CLIP which efficiently learns visual concepts from natural language supervision. CLIP can be applied to any visual classification benchmark by simply providing the names of the visual categories to be recognized, similar to the "zero-shot" capabilities of GPT-2 and GPT-3.

January 5, 2021
15 minute read

- ## CLIP: Connecting Text and Images

**1. Contrastive pre-training**



**2. Create dataset classifier from label text**



**3. Use for zero-shot prediction**

- **CLIP: Image-Text Match**



Cosine similarity between text and image features

# History of Computer Vision

DALL-E
Creating Images from Text

- **Low-level Vision**

# SenseTime – Pioneer in Deep Learning and Computer Vision

**Founding** of CUHK MMLab by Prof. Tang Xiao'ou

CUHK MMLab is the **earliest Chinese** research team to work on **deep learning**

MMLab's CV Papers published **almost equal to** the total amount published by all Chinese universities

MMLab's deep learning papers published in CV conferences are **almost half of** the total amount papers published **globally**

MMLab's DeepID algorithm **exceeds human eye accuracy** for the first time in history

SenseTime became **1st** Chinese company to **win** ImageNet competition, and was ranked top in video analysis

IMAGENET

SenseTime incorporated

SenseTime **surpassed Google and Facebook** in CVPR & ICCV [1] Paper Submission

SenseTime **won three** world champions in five key ImageNet competitions

IMAGENET

SenseTime and its joint labs published 62 papers and 57 papers in 2019 CVPR and ICCV respectively, ranking 1st in the world [2]

**MIT- SenseTime Alliance** was formed to support AI Research

SenseTime and its joint labs published 62 papers in CVPR 2020

| 2001 | 2004-2008 | 2011 | 2011-2013 | 2014.6 | 2014.10 | 2015.11 | 2016.9 | 2017.10 | 2018.2 | 2019.7 | 2020.6 |

| 2011.8 | 2012.12 | 2013.3 | 2013.12 | 2014.11 | 2016.3 | 2016.7 | 2017.3 | 2018 | 2019.3 | 2020 |

Microsoft significantly improved the accuracy of Deep Learning-driven voice recognition

Hinton, the originator of Deep Learning, won the ImageNet visual recognition competition

Google employed Prof. Hinton for US$50 million

Facebook established an AI Lab in New York and employed Yann LeCun

Google acquired DeepMind, a Deep Learning company, for US$660 million

DeepMind's AlphaGo, the Go AI, beats Lee Sedol in Go

GM acquired Cruise Automation, a startup in autonomous vehicle technology for US$1 billion

Softbank launched a US$100 billion investment fund that focuses on AI and completed US$32 billion acquisition of ARM

Intel acquired Mobileye, a leader in computer vision for autonomous driving technology

Tesla launched self-driving vehicle. Waymo launched self-driving taxi service

NVIDIA launches its own GPU cloud

Hinton, LeCun & Bengio received the Turing Award, the "Nobel Prize" of Computing

Global tech companies are exploring the use of AI to combat the outbreak of COVID-19.

(1) CVPR、ICCV、ECCV are the top 3 computer vision conferences worldwide with highest impact factor They accept the best work on computer vision and deep learning

(2) Based on statistics released by different companies and organizations to date

# How to Generate the Best AI

Fundamental research & technological capabilities determine rate of innovation

## Expertise

Large amount of high quality data fuels the algorithm iteration

## Data

Super fast computing power ensures speed of training

## Computing Power

Vertical partnerships ensure technology and data feedback for adaptive improvement

## Positive Feedback Loop

**SenseTime Excels at All of These Core Capabilities**

# SenseTime – World Leading AI Innovation Platform

## Smart City

**Smart Surveillance**

**Smart City Management System**

**Smart Traffic Management**

**Fire Detection**

**Smart Crowd Management**

**Abnormal Behavior Detection**

**Garbage Detection**
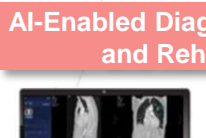
**Illegal Parking Detection**

**Illegal Occupation Detection**

**Abnormal Objects Detection on Road**

## Business Intelligence

**Retail Analytics Solutions**

**Intelligent Hotel Check in System**

**Smart Airport Solution**

**Smart Metro Solution**

**Smart Office Management System**

**Smart Tourism Area Management**

**Smart Entertainment Solution**

**Smart Campus Solution**

**Smart Amusement Park Solution**

**Real Estate Sales Management**

## Mobile Solution

**Face Unlock**

**Photo Processing**

**Image Super Resolution**

**3D Face Beautification**

### AR Platform

**AR Live Streaming**

**AR Game**

**AR Classroom**

**AR Effect**

## Autonomous Driving

**Guide Line Prediction**

**Human Face Prediction**

**Lane Detection**

**Front Vehicle Detection**

### Intelligence Cabin Sensing

**Face Unlock**

**Gaze Tracking**

**Gesture Tracking**

**Drowsiness detection**

## AI Education Package

**AI Textbook**

**AI Experiment Platform**

**AI RobotCar**

**AI Lab**

### Remote Sensing

**Road Network Extraction**

**Cloud and Snow Detection**

### AI-Enabled Diagnosis, Treatment and Rehabilitation

**Lung AI Application**

**Pathology Application**

WONG KAR-WAI'S

## IN THE MOOD
## FOR LOVE