



Introduction to Computer Vision

Lecture 1 - Overview

Prof. He Wang

About Me

- 王鹤
- Assistant Professor in Center on Frontiers of Computing Studies (CFCS)
- Joined PKU in September, 2021
- Received Ph.D. from Stanford in 2021
- Received Bachelor from Tsinghua in 2014
- Our lab: *Embodied Perception and Interaction (EPIC) Lab*
- Research interest: 3D vision, Robotics
- Homepage: <https://hughw19.github.io/>



北京大学前沿计算研究中心
Center on Frontiers of Computing Studies, Peking University



Course Logistics

Objective: A Great Course on Computer Vision

- A self-included beginning course on computer vision
- A broad coverage of classic and deep vision from a modern perspective, to distinguish from online available vision courses
- To lay a solid foundation for applying and researching in computer vision

Logistics

- Instructor
 - He Wang (hewang@pku.edu.cn)
 - Office Hour: Friday 5:00PM - 6:00PM or under appoint.
 - Office location: Room 106-1, Courtyard No.5, Jingyuan
- TAs:
 - Mi Yan (dorisyan@pku.edu.cn)
 - Hao Shen (2301112029@pku.edu.cn)
 - Yuxing Chen (yuxingc_20@stu.pku.edu.cn)
- Class Time & Location
 - Wednesday 3:10PM - 6:00PM
 - Room 507, Teaching Building 2, Peking University

Prerequisite

- Math:
 - Calculus
 - Linear Algebra
 - Basics of probability and statistics
- Proficiency in Python
- Optionally, have taken *Introduction to AI* or know some basic knowledge of machine learning and neural networks.

Books and References

- No required textbooks.
- Books for references:
 - On Deep Learning:
 - Ian Goodfellow and Yoshua Bengio and Aaron Courville. *Deep Learning*. MIT Press, 2016.
 - On Classic Computer Vision:
 - R. Szeliski. *Computer Vision: Algorithms and Applications*. Springer, 2011.
 - D. A. Forsyth and J. Ponce. *Computer Vision: A Modern Approach (2nd Edition)*. Prentice Hall, 2011.

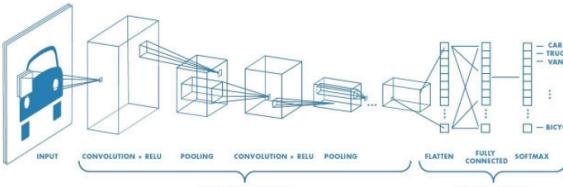
Books and References

- The most effective way to learn and check things:
 - Just google it!
 - Search in wikipedia!

Published in Towards Data Science · Follow

Sumit Saha
Dec 16, 2018 · 7 min read

A Comprehensive Guide to Convolutional Neural Networks — the ELI5 way



Artificial Intelligence has been witnessing a monumental growth in bridging the gap between the capabilities of humans and machines. Researchers and enthusiasts alike, work on numerous aspects of the field to make amazing things happen. One of many such areas is the domain of Computer Vision.

The agenda for this field is to enable machines to view the world as humans do, perceive it in a similar manner and even use the knowledge for a multitude of tasks such as Image & Video recognition, Image Analysis & Classification, Media Recreation, Recommendation Systems, Natural Language Processing, etc. The advancements in Computer Vision with Deep Learning has been constructed and perfected with time, primarily over one particular algorithm — a Convolutional Neural Network.

Not logged in Talk Contributions Create account Log in

Read Edit View history Search Wikipedia

Generative adversarial network

From Wikipedia, the free encyclopedia

Not to be confused with Adversarial machine learning.

A generative adversarial network (GAN) is a class of machine learning frameworks designed by Ian Goodfellow and his colleagues in June 2014.^[1] Two neural networks contest with each other in a game (in the form of a zero-sum game, where one agent's gain is another agent's loss).

Given a training set, this technique learns to generate new data with the same statistics as the training set. For example, a GAN trained on photographs can generate new photographs that look at least superficially authentic to human observers, having many realistic characteristics. Though originally proposed as a form of generative model for unsupervised learning, GANs have also proved useful for semi-supervised learning,^[2] fully supervised learning,^[3] and reinforcement learning.^[4]

The core idea of a GAN is based on the "indirect" training through the discriminator, another neural network that is able to tell how much an input is "realistic", which itself is also being updated dynamically.^[5] This basically means that the generator is not trained to minimize the distance to a specific image, but rather to fool the discriminator. This enables the model to learn in an unsupervised manner.

Part of a series on **Machine learning and data mining**

- Problems
- Supervised learning (classification + regression)
- Clustering
- Dimensionality reduction
- Structured prediction
- Anomaly detection
- Artificial neural network
- Autoencoder - Cognitive computing
- Deep learning - DeepDream
- Multilayer perceptron - RNN (LSTM - GRU)
- ESN - Restricted Boltzmann machine - GAN
- SOM - Convolutional neural network (U-Net)
- Transformer (Vision) - Spiking neural network
- Memristor - Electromechanical RAM (ECRAM)

Reinforcement learning

Theory

Machine-learning venues

Related articles

Method [edit]

The *generative network* generates candidates while the *discriminative network* evaluates them.^[1] The contest operates in terms of data distributions. Typically, the generative network learns to map from a latent space to a data distribution of interest, while the discriminative network distinguishes candidates produced by the generator from the true data distribution. The generative network's training objective is to increase the error rate of the discriminative network (i.e., "fool" the discriminator network by producing novel candidates that the discriminator thinks are not synthesized (are part of the true data distribution)).^{[1][6]}

A known dataset serves as the initial training data for the discriminator. Training it involves presenting it with samples from the training dataset, until it achieves acceptable accuracy. The generator trains based on whether it succeeds in fooling the discriminator. Typically the generator is seeded with

Courseworks and Grading Policy

- 4 assignments: each 10%, in total **40%**
- 1 midterm exam: **30%**
- 1 final exam: **30%**
- Class/discussion board participation: up to **5% bonus**

- Look up class schedule ([https://pku-epic.github.io/
Intro2CV_2024/schedule/](https://pku-epic.github.io/Intro2CV_2024/schedule/)) for release and due dates.

Assignments

- Late policy for assignments
 - If 1 day (0 - 24 hours) past the deadline, 15% off
 - If 2 day (24 - 48 hours) past the deadline, 30% off
 - Zero credit if more than 2 days.

Midterm and Final Exams

- Midterm exam will be held in class.
- Final exam will be held in the afternoon of June 19.
- 1-page cheat sheet is allowed for both.

Course Website

- Website accessible to everyone: [https://pku-epic.github.io/
Intro2CV_2024/schedule/](https://pku-epic.github.io/Intro2CV_2024/schedule/)
- Internal course website in <https://course.pku.edu.cn/>
 - Slides/videos download
 - Assignment submission
 - Grades
 - Discussion board

My Experience back to Stanford

The screenshot shows a Piazza discussion board for the CS 231N course. The top navigation bar includes links for LIVE Q&A, Drafts, google_cloud, midterm, project, other, hw1, lectures, office_hour, hw2, hw3, pytorch, tensorflow, hyperquest, and He Wang. A search bar and a user profile for He Wang are also present.

The main content area displays a list of posts filtered by 'hw1'. A banner at the top states: 'This class has been made inactive. No posts will be allowed until an instructor reactivates the class.' The first post is titled '[HW1 features] Neural Networks with Image features' and asks about achieving 55% or 60% credit. It includes a statement from the instructor: 'you should easily be able to achieve over 55% classification accuracy on the test set; our best model achieves about 60% classification accuracy.' Below this is a section for 'the students' answer' and another for 'the instructors' answer', both of which are currently empty.

Further down, there is a 'followup discussions' section. A post by Rishab Mehra asks if achieving 55% on the test set is sufficient. Aman Peddada responds that more than 55% is better. Rishab Mehra then explains that the test set needs to be 'tuned' to achieve over 55%, which seems counter-intuitive given the definition of a test set.

The bottom of the page shows a form to 'Start a new followup discussion'.

Key posts visible include:

- [HW1 features] Neural Networks with Image features (4/24/17)
- [HW1 softmax] 0.5 coefficient for regul... (4/22/17)
- Time taken for HW1 (4/22/17)
- how to submit HW1? (4/21/17)
- Gradient of bias (4/21/17)
- [HW1 features] Neural Networks with I... (4/21/17)
- [HW1] ipython notebook stuck.. after i... (4/21/17)
- [HW1 NN] 2nn, without tuning, accurac... (4/21/17)
- Are we allowed to have additional .py fi... (4/21/17)
- 1HW1 NNT Above 48% on validation but (4/21/17)

Discussion Board

- Course discussion board: **share your questions!**
- Will have a constant forum for discussing course material and three assignment forums each for one assignment.
- Your active participation in classes, discussions and my office hours will all count towards bonus.

讨论区	描述	帖子总数
<input type="checkbox"/> 论坛		
<input type="checkbox"/> 课程问题讨论		12
<input type="checkbox"/> 作业一		107
<input type="checkbox"/> 作业二		29
<input type="checkbox"/> 发帖规范	1. 如果是作业相关的帖子, 请在对应论坛内将标题设为具体的小问名称+简要概括内容, 如 "Task 1(a) 题干描述有误", 0 便于其他同学查找 2. 如果是询问某次作业中一个函数是否可以使用, 请在对应论坛内助教发的帖子下回复具体小问+函数名称, 如 "Task 1(a) np.vectorize" 3. 请大家发帖前先查看是否已有过相同或者相似的问题, 如果是相似问题请直接在已有帖子下回复。	0
<input type="checkbox"/> 期中线上答疑	期中考范围是Lec1~Lec8。英文题面, 作答中英文皆可。题型: 多选, 判断, 简答, 计算。	17
<input type="checkbox"/> 作业三		40
<input type="checkbox"/> 作业四		53
<input type="checkbox"/> 期末答疑		10

WeChat Group

- WeChat group:
 - Use for notifications and announcements
 - Not an ideal place for tracking each individual questions.
 - Not recommended to WeChat me or TAs in person to ask questions that may also be interesting the other students.
Use discussion board!



群聊: 2024春计算机视
觉导论



该二维码7天内(2月28日前)有效，重新进入将更新

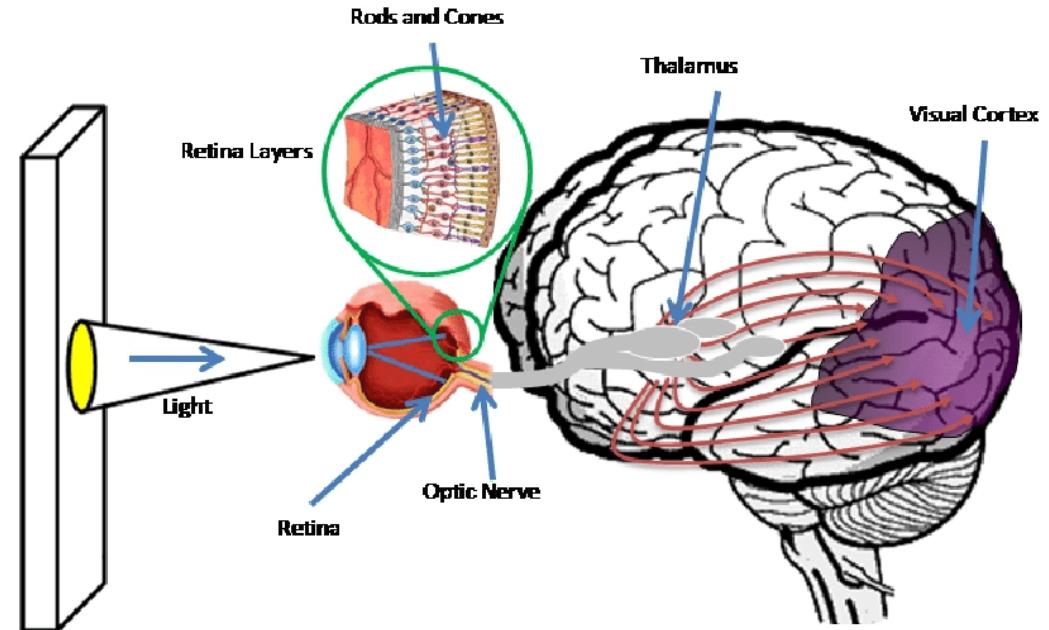
What is Vision?

What is Vision?



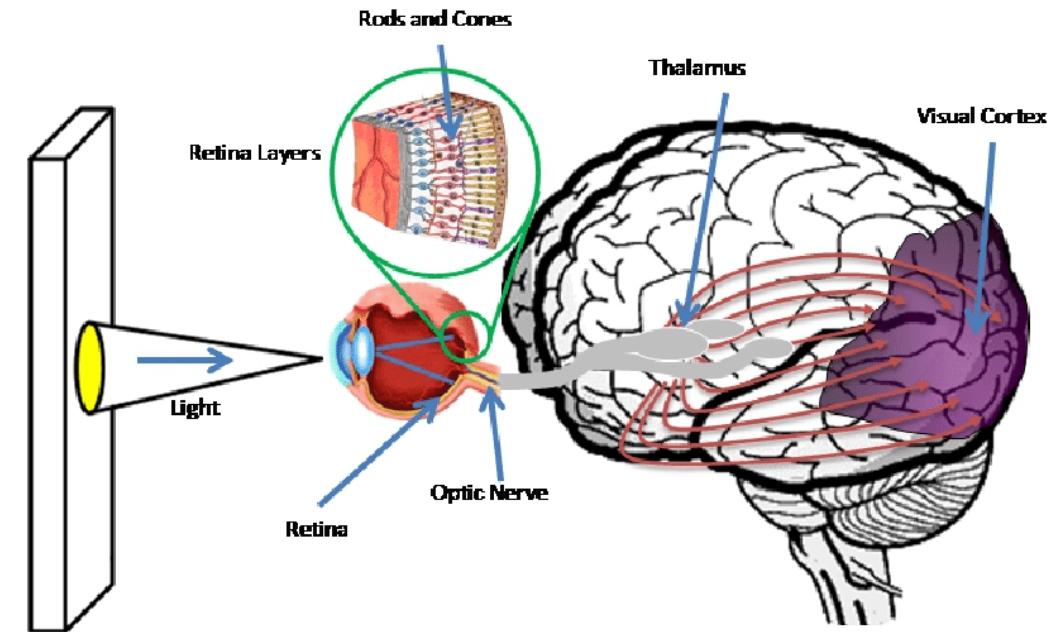
Human Visual System

- The visual system comprises
 - Eyes (sensory organ)
 - Parts of the central nervous system
 - Retina layers
 - Optic nerve
 - Optic tract
 - Visual cortex



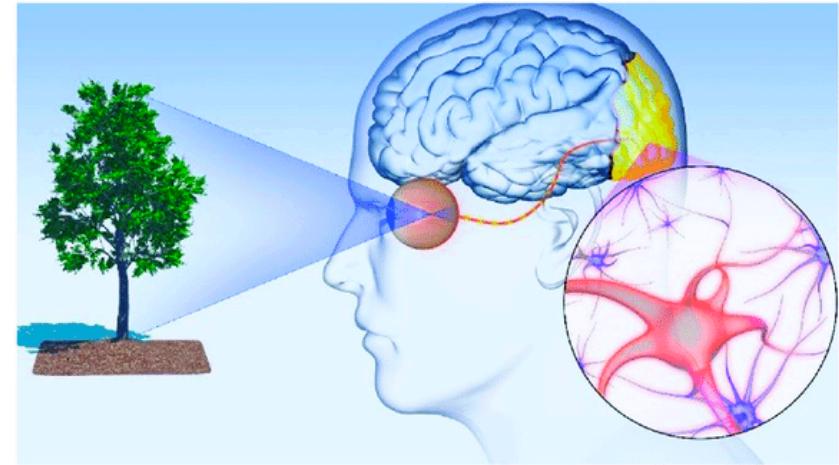
Human Visual System

- The visual system comprises
 - Eyes (sensory organ)
 - Parts of the central nervous system
 - Retina layers
 - Optic nerve
 - Optic tract
 - Visual cortex
- Visual pathway: visual field → retina → optic nerve → ... → optic tract → ... → visual cortex



Human Visual System

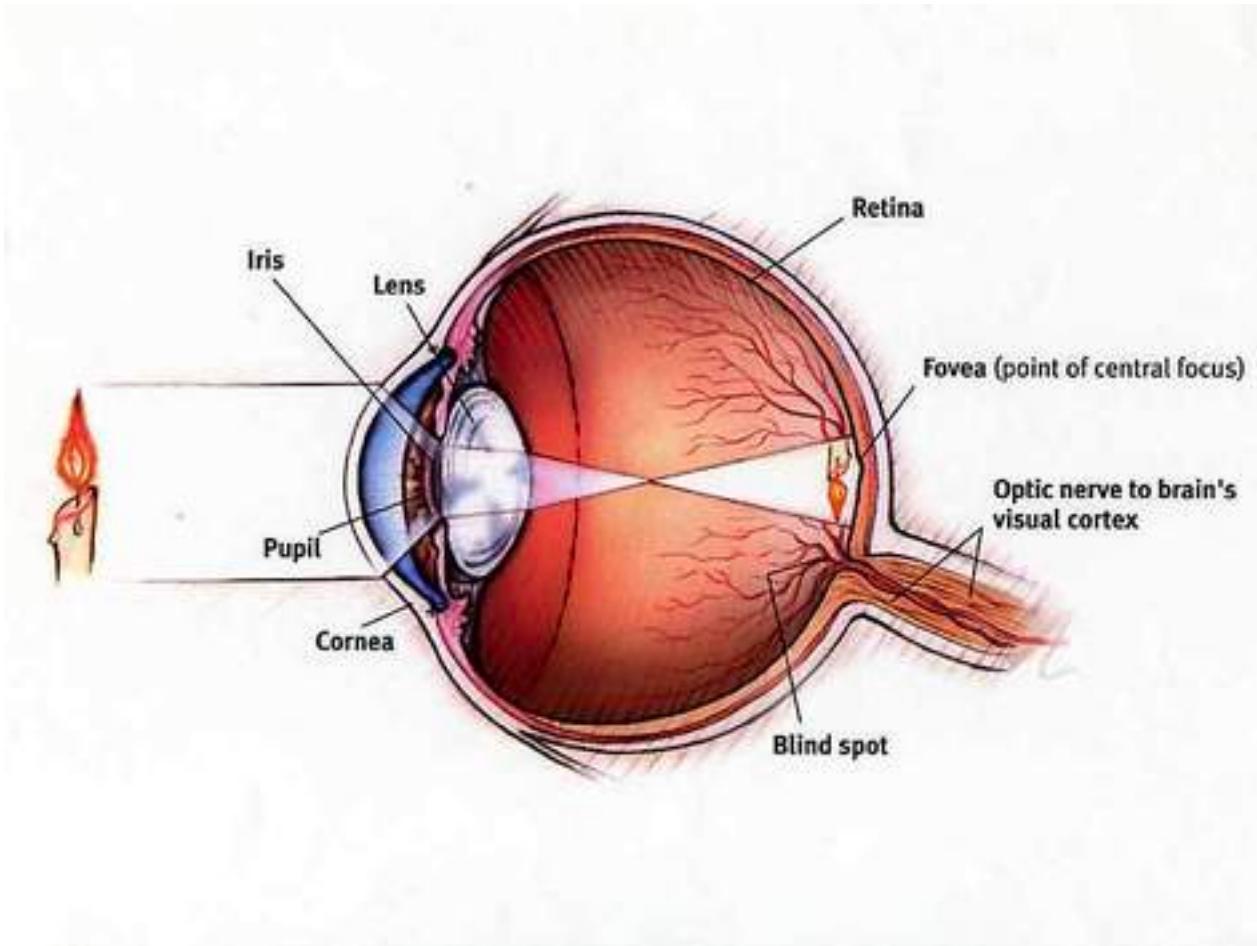
- 83% information comes from vision (11% from audio, the others from smell, touch and taste).
- Carry outs a number of complex tasks:
 - visual sensation
 - visual perception
 - and visual motor coordination.



Our Visual System Needs to Do

- the reception of light
 - the formation of monocular neural representations
 - color vision
 - stereopsis and assessment of distances to and between objects
- }
- Visual sensation
-
- pattern recognition
 - the identification of particular object of interest
 - motion perceptions
 - the analysis and integration of visual information
 - accurate motor coordination under visual guidance
 - ...
- }
- Visual perception
-
- Integrating proprioception and vision signals (eye sensory feedbacks)
 - Moving body parts as required to accomplish intended actions, e.g. eye-hand coordination in writing, eye-muscle coordination in sports
 - ...
- }
- Visual motor coordination

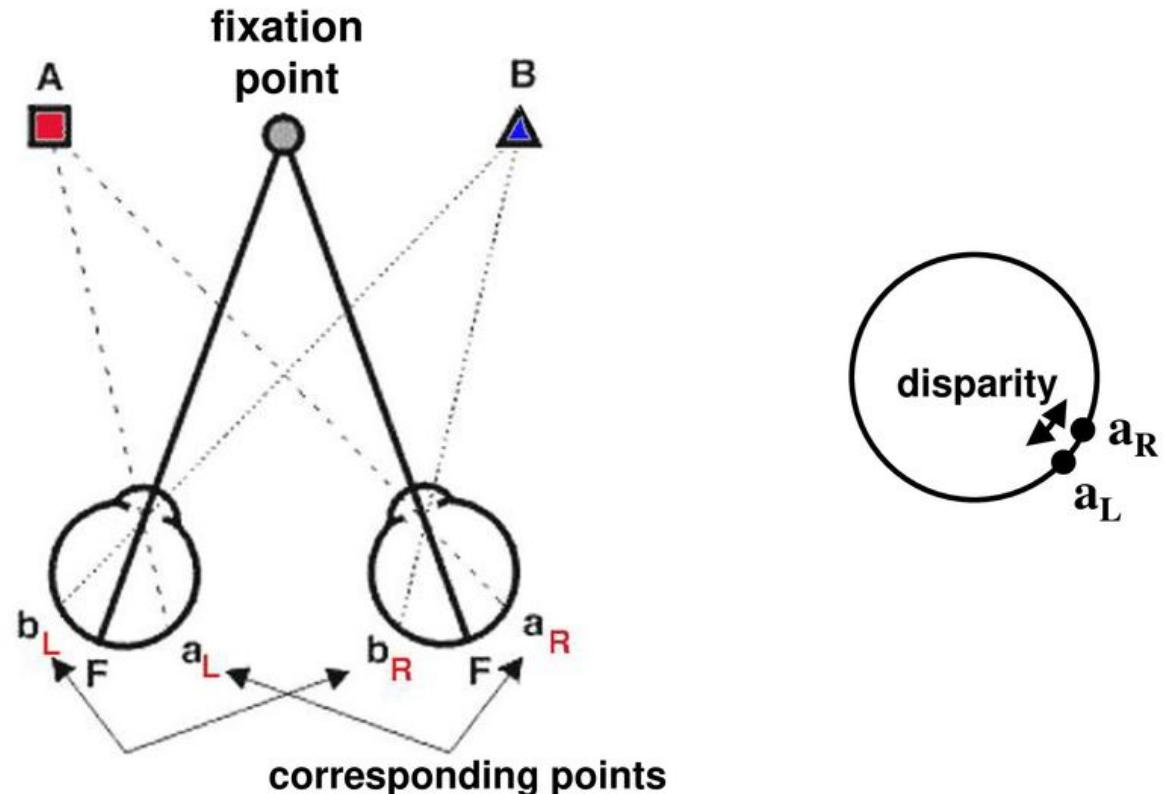
Visual Sensation: Monocular Vision



Visual Sensation: Binocular Vision and Stereopsis

- Human eyes are binocular.
- Senses distances through stereopsis

Human stereo geometry



http://webvision.med.utah.edu/space_perception.html

S. Birchfield, Clemson Univ., ECE 847, <http://www.ces.clemson.edu/~stb/ece847>

Visual Perception

- Definition of visual perception in *Vision Science*:
 - the process of acquiring knowledge about environmental objects and events by extracting information from the light they emit or reflect.



Visual Perception

- Concerns the acquisition of knowledge.
 - Fundamentally a cognitive activity.
 - Distinct from purely optical processes such as photographic ones.
 - Vision = eyes = camera? No, cameras have no perceptual capabilities at all.



Visual Perception

- Concerns the acquisition of knowledge.
- The knowledge achieved by visual perception concerns objects and events in the environment.



Examples of Perception: Motion Perception

- Motion perception: the ability of the nervous system to discern the distance and speed of a moving object in relation to the eye that is seeing the object.



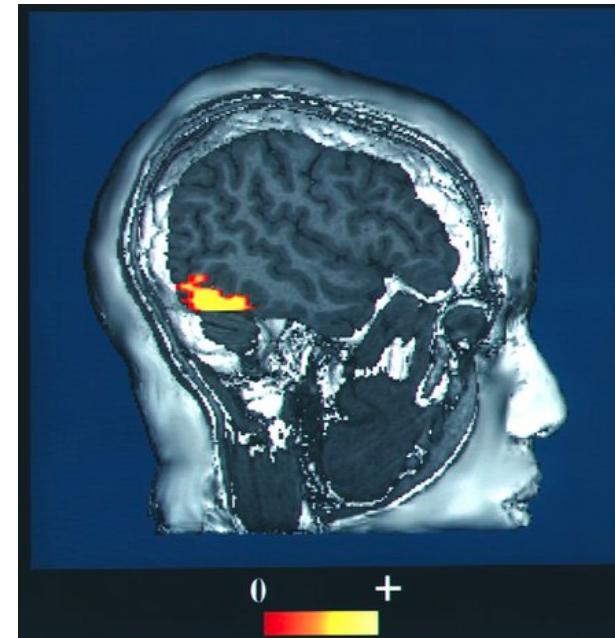
Wikipedia contributors. "Visual system." *Wikipedia, The Free Encyclopedia*. Wikipedia, The Free Encyclopedia, 18 Jan. 2022. Web. 22 Feb. 2022.
Hülsdünker, Thorben, Martin Ostermann, and Andreas Mierau. "The speed of neural visual motion perception and processing determines the visuomotor reaction time of young elite table tennis athletes." *Frontiers in behavioral neuroscience* 13 (2019): 165.

Examples of Perception: Pattern Recognition

- One great example of pattern recognition is facial recognition.



Facial recognition: detect faces,
perceive emotion, distinguish similar faces
(Dimitri Otis | Getty Images)



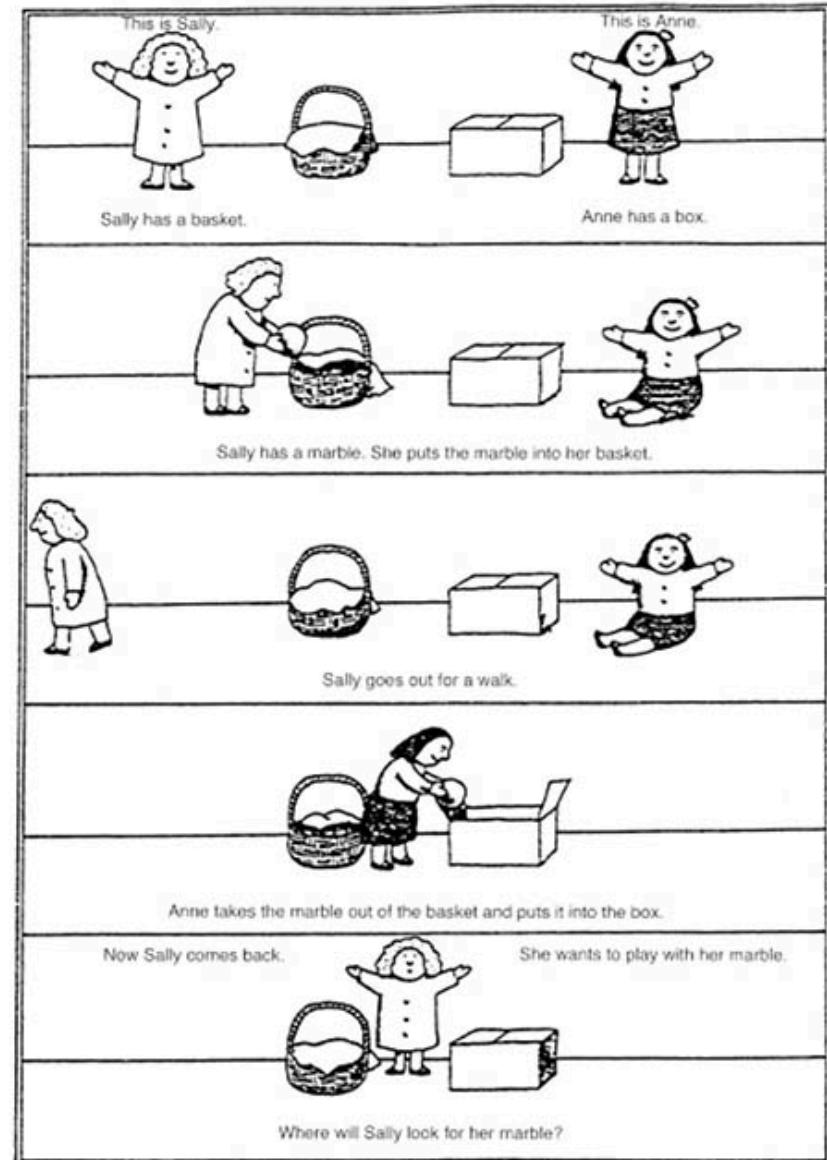
Corresponding area in the brain

Examples of Perception: Visual Cognition

- One example of visual cognition is false belief.

- Sally-Anne Test

Sally takes a marble and hides it in her basket. She then "leaves" the room and goes for a walk. While she is away, Anne takes the marble out of Sally's basket and puts it in her own box. Sally is then reintroduced and the child is asked the key question, the *Belief Question*: "Where will Sally look for her marble?"



Wikipedia contributors. "Sally–Anne test." Wikipedia, The Free Encyclopedia. Wikipedia, The Free Encyclopedia, 9 Feb. 2022. Web. 22 Feb. 2022.

Kosinski, Michal. "Theory of Mind May Have Spontaneously Emerged in Large Language Models." arXiv preprint arXiv:2302.02083 (2023).

Examples of Perception: Visual Cognition

- One example of visual cognition is false belief.

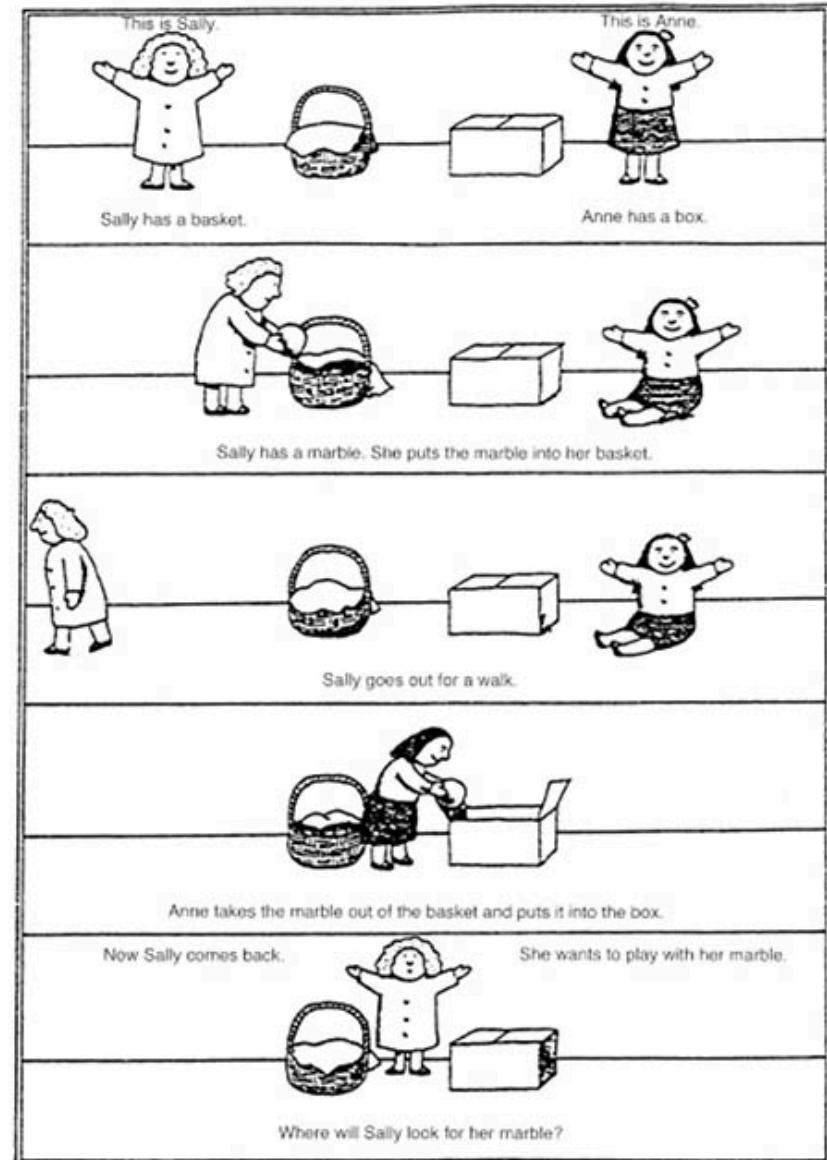
- Sally-Anne Test

Sally takes a marble and hides it in her basket. She then "leaves" the room and goes for a walk. While she is away, Anne takes the marble out of Sally's basket and puts it in her own box. Sally is then reintroduced and the child is asked the key question, the *Belief Question*: "Where will Sally look for her marble?"

- Does ChatGPT have theory of mind?

Wikipedia contributors. "Sally–Anne test." Wikipedia, The Free Encyclopedia. Wikipedia, The Free Encyclopedia, 9 Feb. 2022. Web. 22 Feb. 2022.

Kosinski, Michal. "Theory of Mind May Have Spontaneously Emerged in Large Language Models." arXiv preprint arXiv:2302.02083 (2023).



Visual Motor Coordination/Integration

- Visual motor control is the ability to coordinate visual information with motor output, where the eyes provide sensory feedback to adjust body motion.
- It is crucial for coordinating the hands, legs, and the rest of the body's movements with what the eyes perceive.

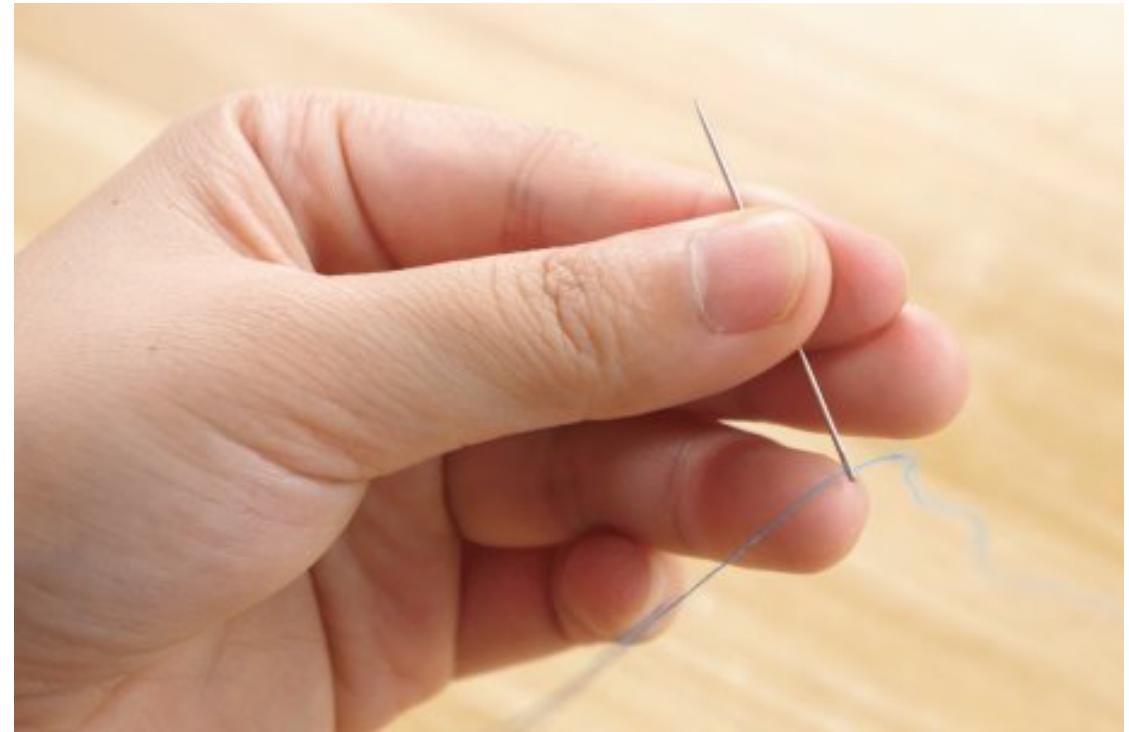
OT Mom's Visual Perception Activities

Why Visual-Motor Integration Is sooooo Important For Handwriting



OT Mom Learning Activities

Examples of Visuomotor coordination: Eye-Hand Coordination



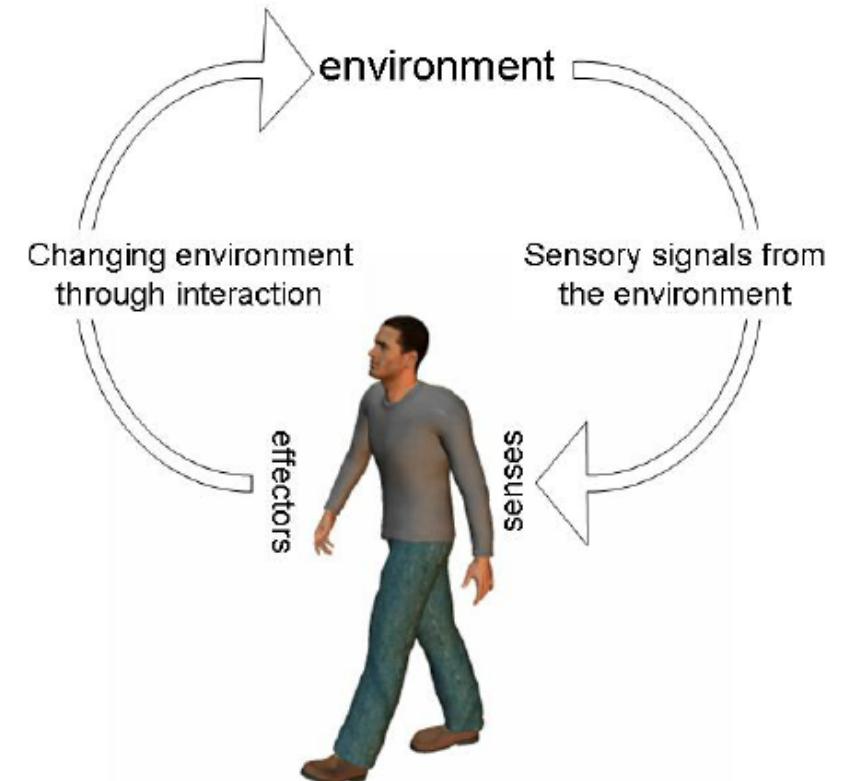
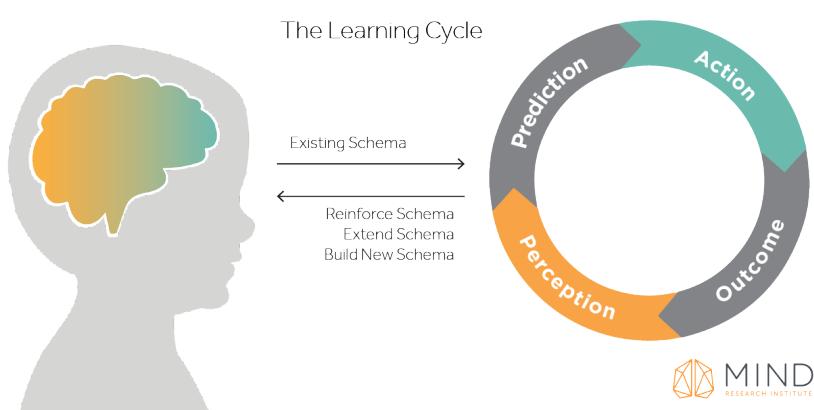
Use vision signal as feedback to perform closed-loop control

More Examples: Running, Balancing, etc.



Perception-Action Loop

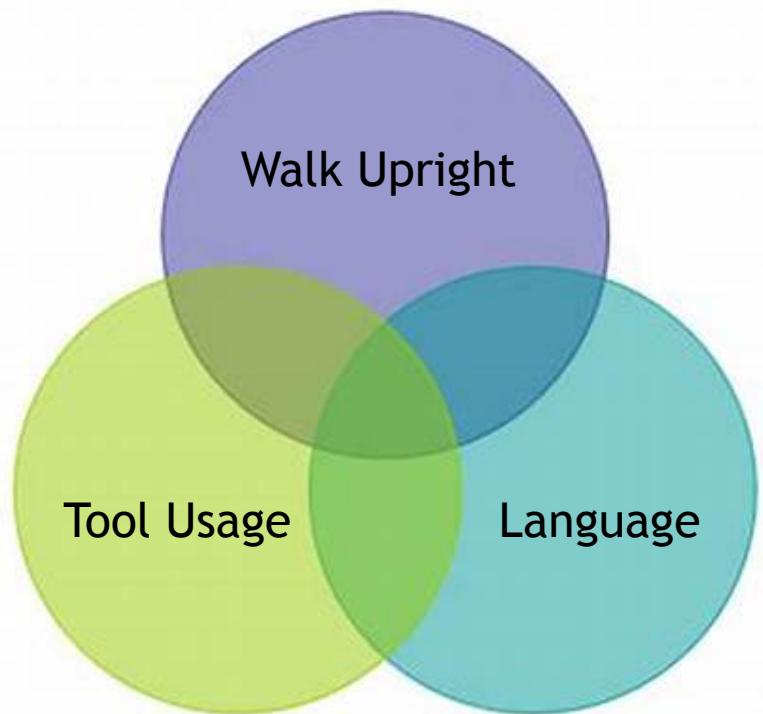
- Perceive, forms hypotheses, and then take action to examine.
- Our brain makes sense of the world around us by creating and testing hypotheses about the way the world works.



Perception-Action Loop

Vision and Language

The keys to evolution of human intelligence:



The interactions between vision and language:

Talk about what you see

Question: what is the nightstand made of ?



1. can't tell it's covered in cloth
2. it appears to be a large red pillow that may be leather
3. I can't tell
4. I can not tell
5. not sure
6. can't tell
7. some kind of metal , it's out of focus
8. Wood
- ...
- 99.0
- 100.I can't see a baggage cart



Grounding according to language description

Expression = "Woman standing in between two guys"



Summary of Human Visual System

- The visual system comprises eyes as sensor organ, which are connected to visual cortex in the brain.
- Visual tasks:
 - sensation
 - processing
 - perception
 - cognition
 - visuomotor coordination
 - interaction with language system
 - and more.

Computer Vision

What is Computer Vision?

- Computer vision deals with
 - acquiring
 - processing and analyzing
 - understanding
 - generating or imagining
- visual data, ...

Visual Data Acquisition



RGB camera



RGB image

Visual Data Acquisition

- Different types of sensors and visual data



RGB camera



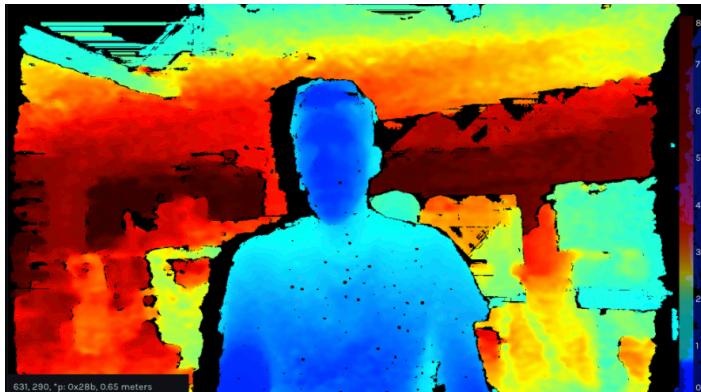
Depth camera



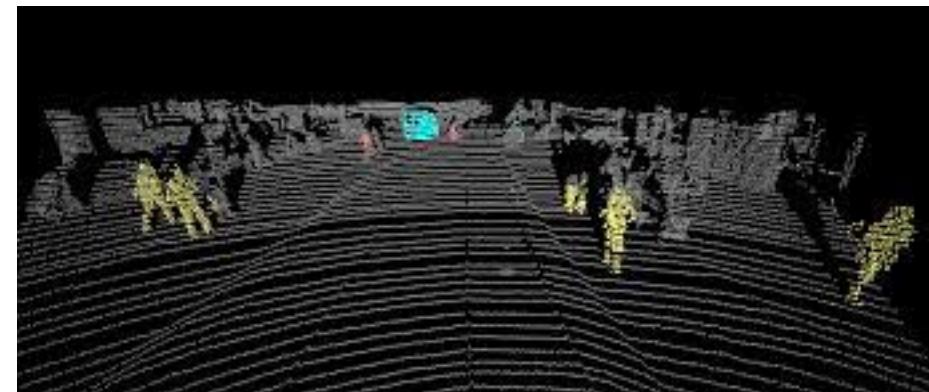
LiDAR



RGB image



Depth image



LiDAR point cloud

Beyond Single Frame and Single View

Stereo
images



Multiview
images

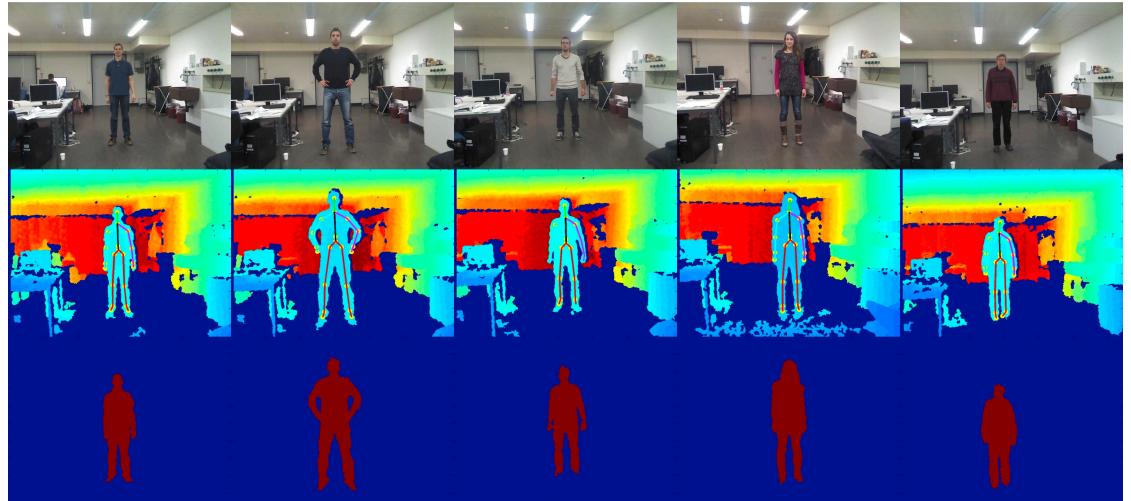


Panoramic images

Beyond Single Frame and Single View

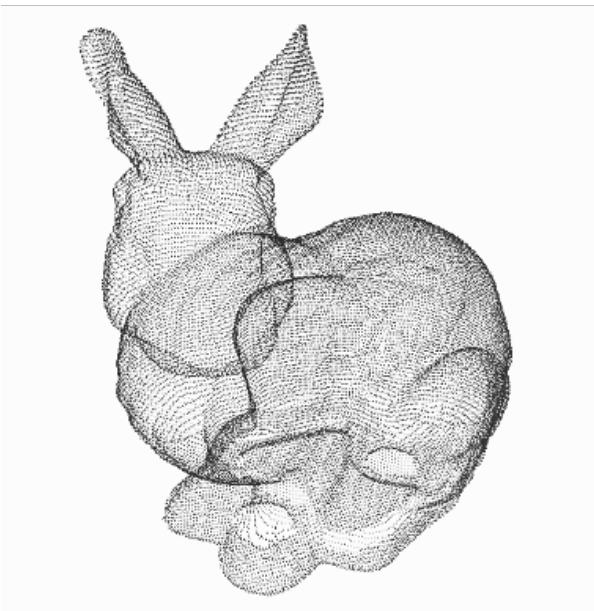


- RGB video

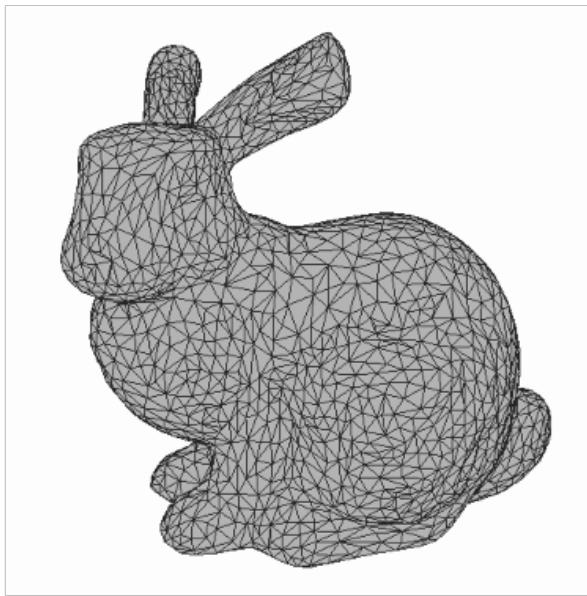


- RGBD video

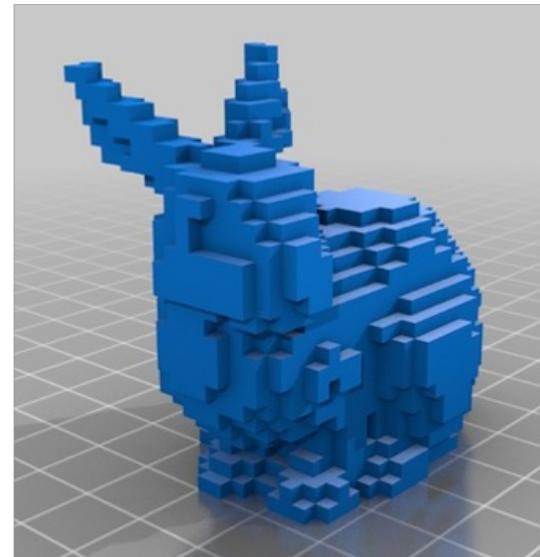
True 3D Visual Data



Point Cloud



Surface Mesh



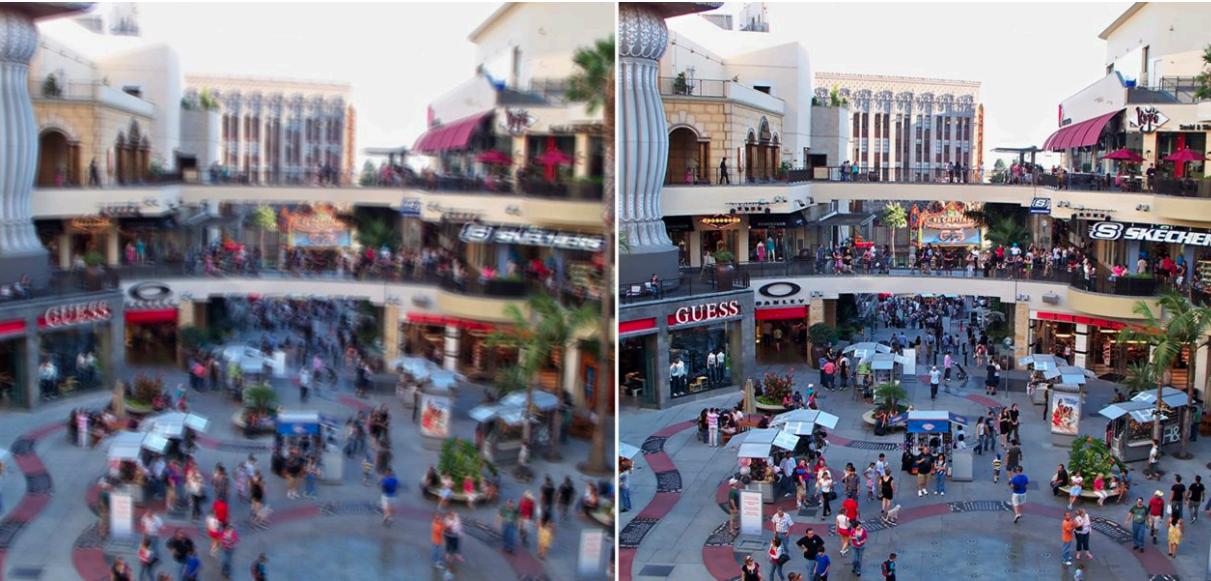
Volumetric

Low-Level Vision: Processing and Feature Extraction

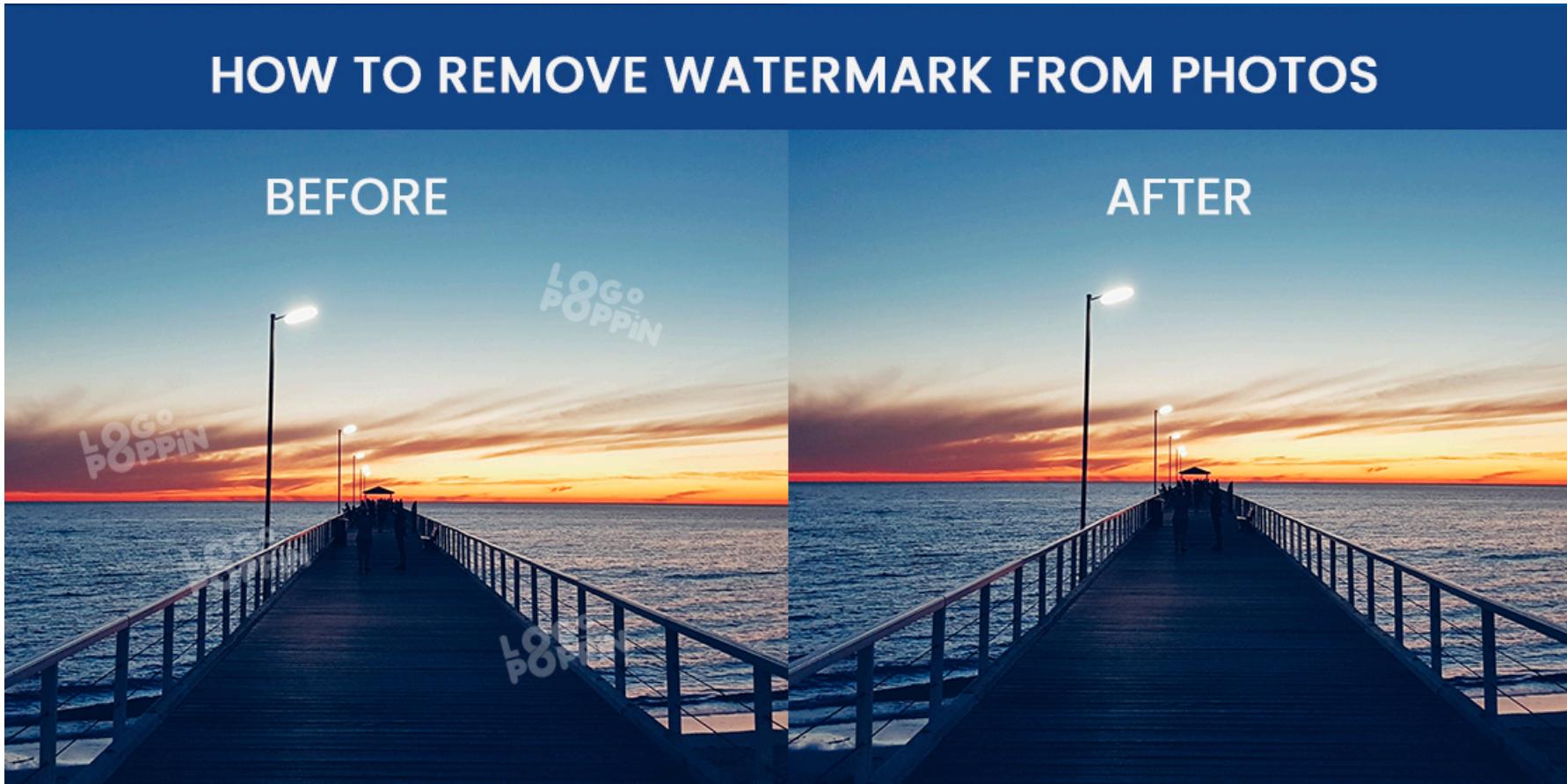
- Low-level vision deals with
 - Image processing
 - image denoising/deblur
 - contrast enhancement
 - ...
 - Feature extractions
 - edge/corner detection
 - optical flow/correspondence



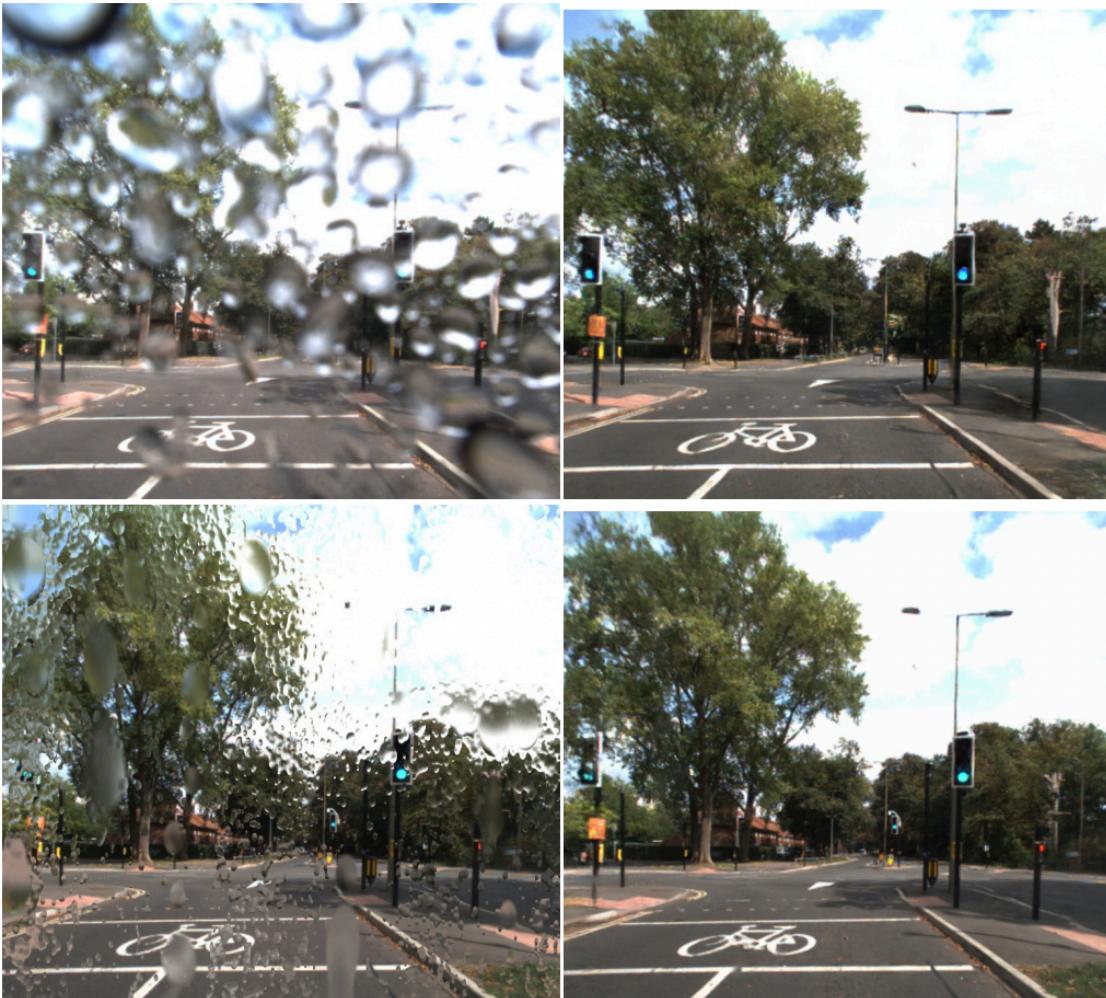
Applications: Motion Deblurring



Applications: Watermark Removal



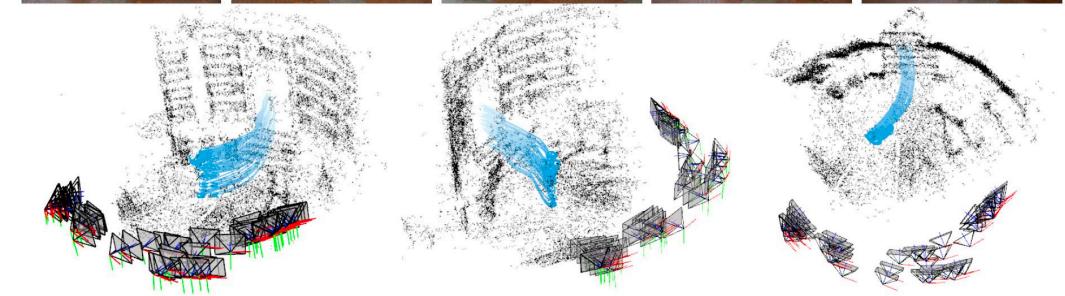
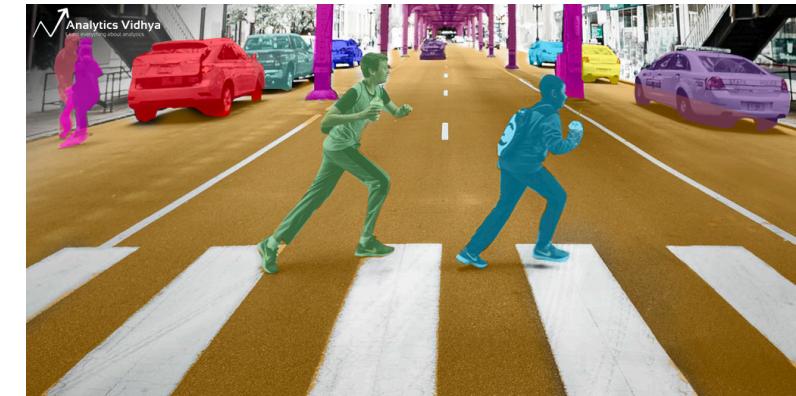
Applications: De-Raining



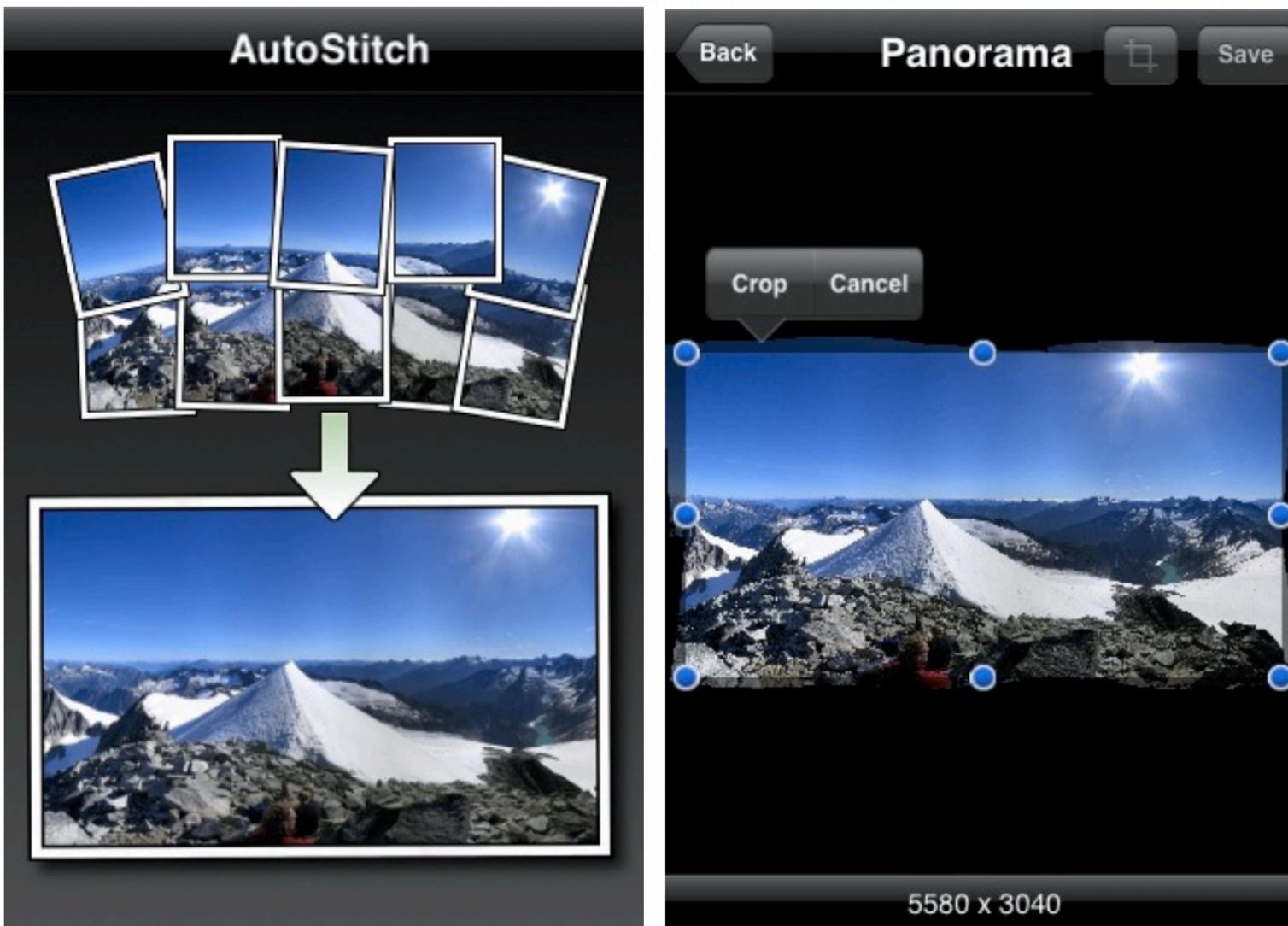
Porav, Horia, Tom Bruls, and Paul Newman. "I can see clearly now: Image restoration via de-raining." *2019 International Conference on Robotics and Automation (ICRA)*. IEEE, 2019.

Mid-Level Vision

- Mid-level vision begins to make inferences about the world based on those measurements.
 - Analyzing local structures (grouping based segmentation, motion analysis, etc.)
 - 3D reconstruction using features obtained from the low-level vision.



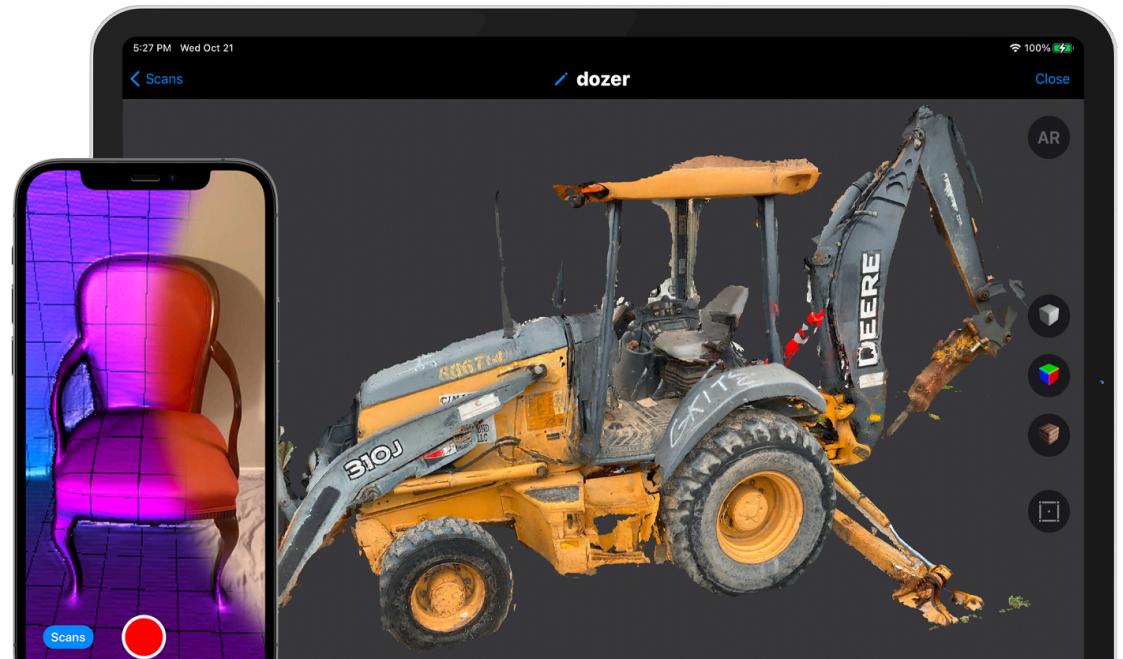
Applications: Panoramic Photography



Applications: 3D Object Scanner



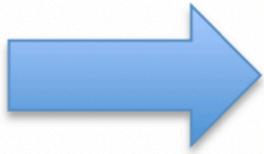
Traditional line scanner



iPad Pro LiDAR Scanner App

Applications: 3D Modeling of Landmarks

From a collection of images, automatically extract features and build a 3D model.

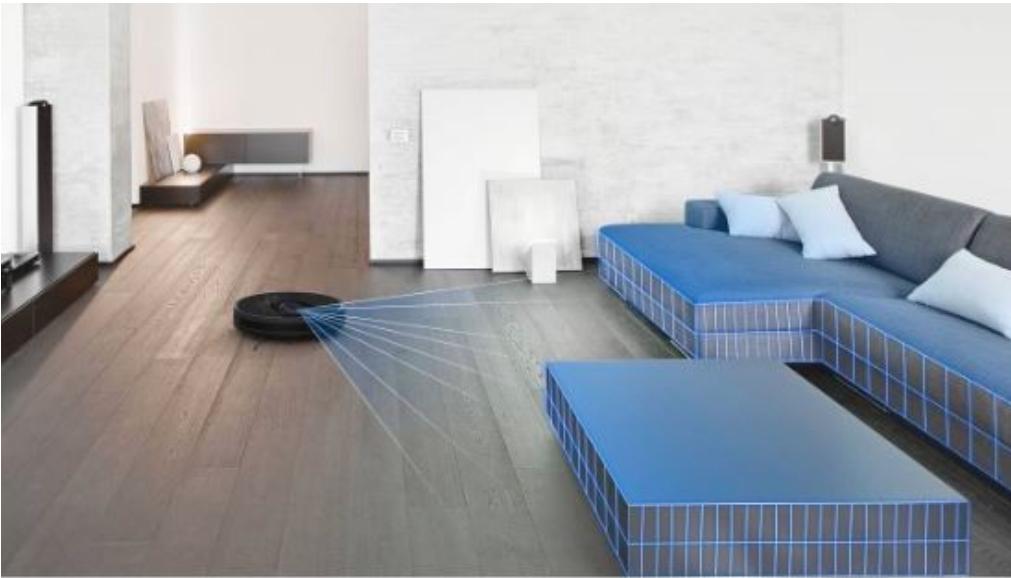


Building Rome in a day.

Agarwal, S., Furukawa, Y., Snavely, N., Simon, I., Curless, B., Seitz, S. M., & Szeliski, R. (2011). Building rome in a day. *Communications of the ACM*, 54(10), 105-112.

Applications: SLAM

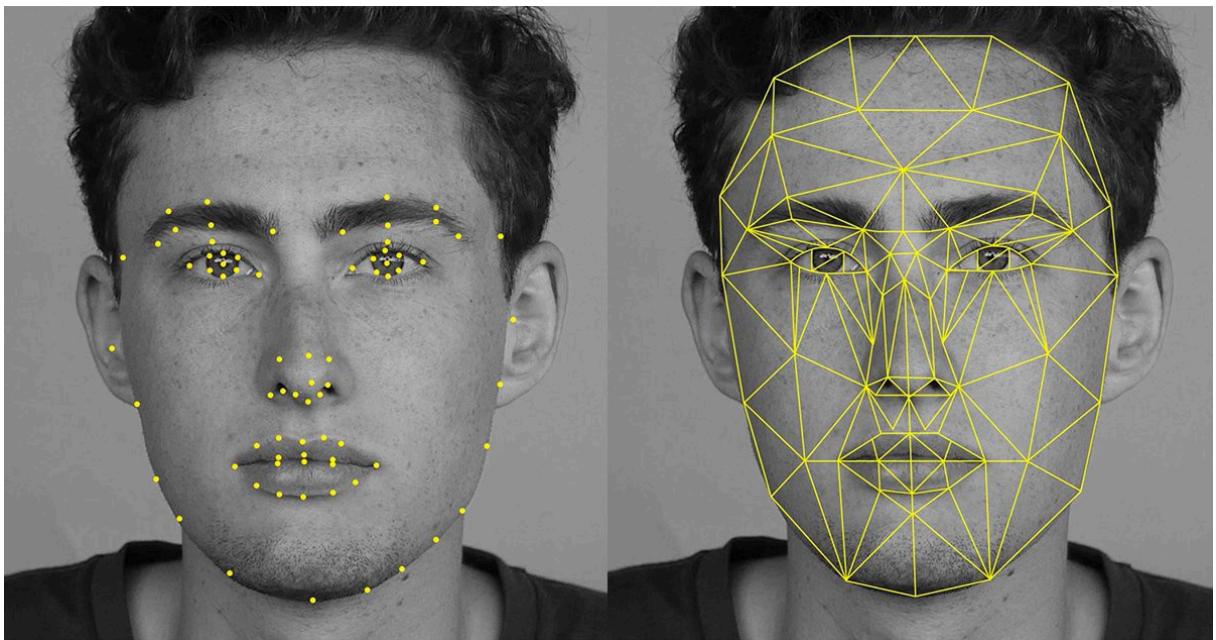
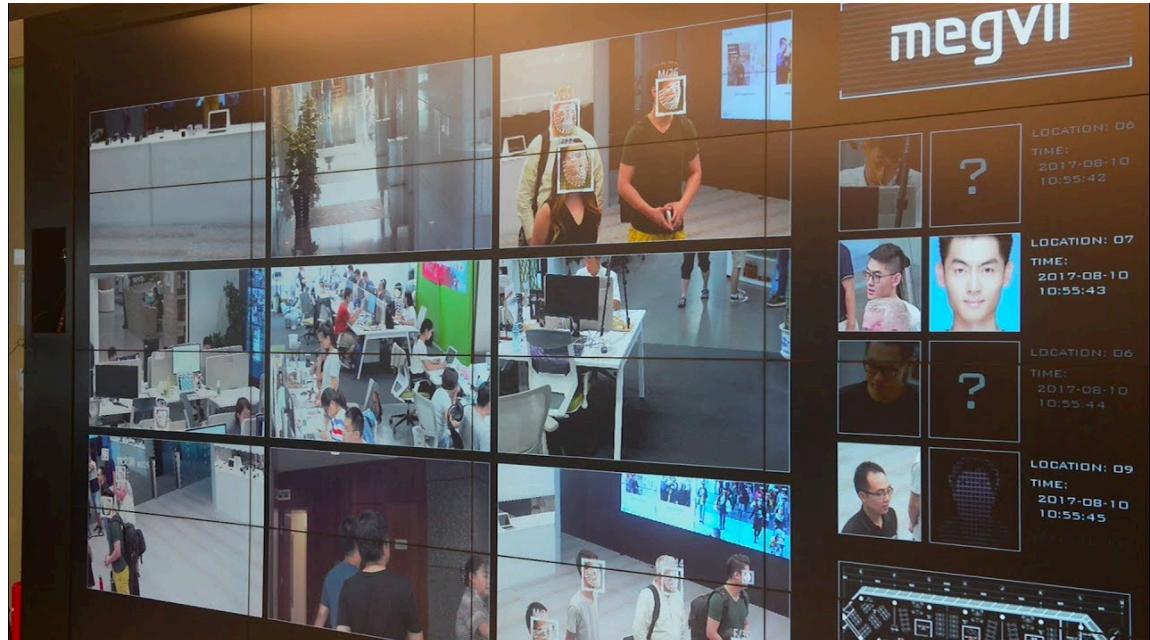
- SLAM = Simultaneous localization and mapping



High-Level Vision: Understanding

- High-level vision analyzes the structure of the external world that produced those images and generates **semantic representation/interpretations**, including
 - object recognition and detection
 - scene understanding
 - activity understanding
 - etc.

Applications: Facial Recognition



Task: Scene Understanding

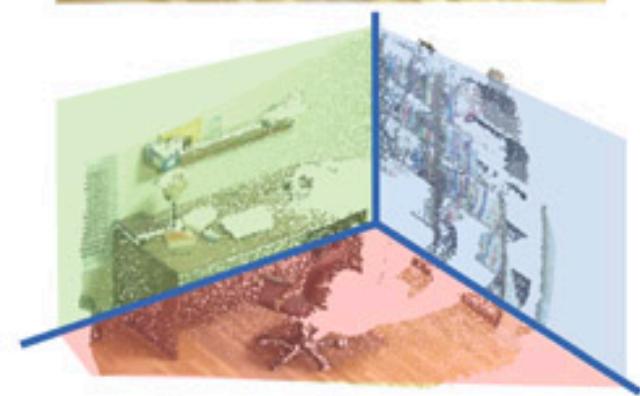
Scene Classification



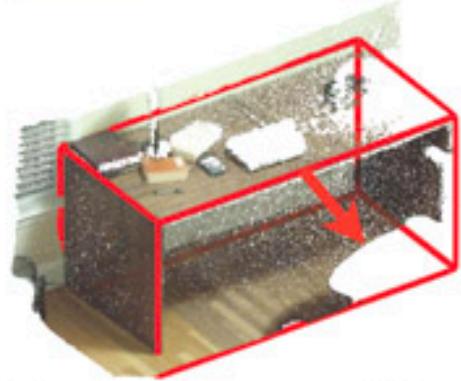
Semantic Segmentation



Room Layout



Detection and Pose



Total Scene Understanding



Applications: Augmented Reality

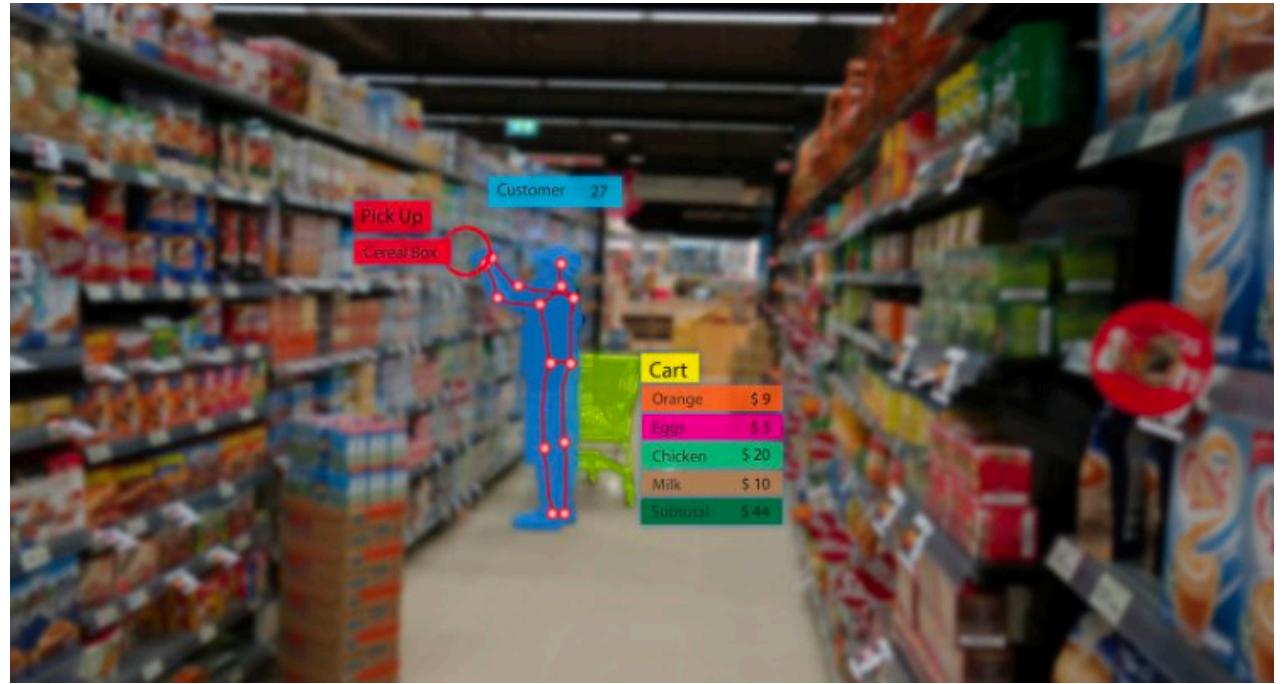


Assisting furniture layout

A photograph of a person's hands holding a white tablet. The screen shows a user interface for a coffee machine, with arrows pointing to different buttons and components. In front of the tablet, a white coffee machine is dispensing coffee into a cup. To the right, there is a red diagonal banner with the text "COMPUTER VISION AND AUGMENTED REALITY". The XRMEET logo is in the top right corner, and contact information is at the bottom right: support@xrmeet.io and www.xrmeet.io.

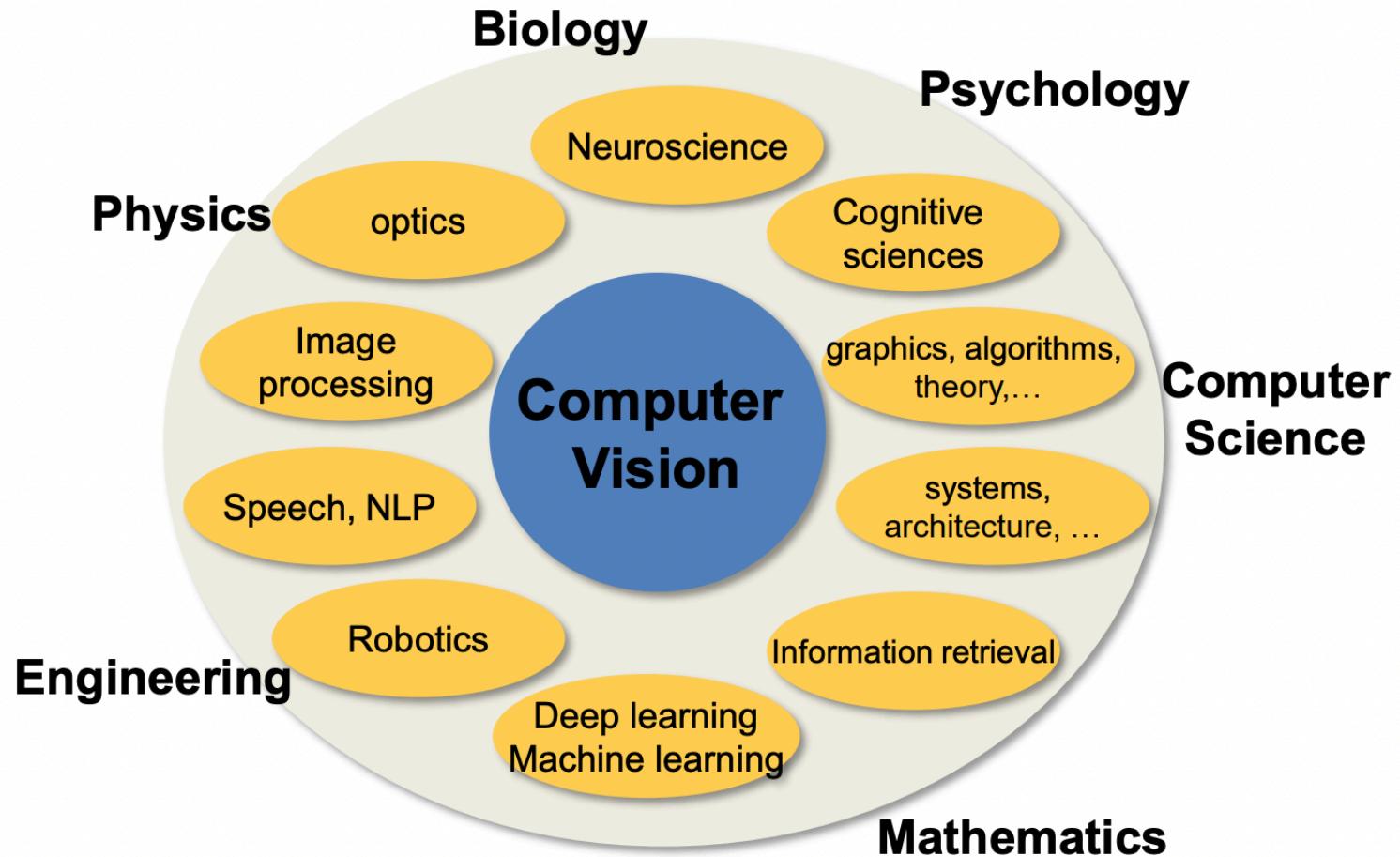
E-Learning

Applications: Amazon Go



- Cashier-free store. A very complicated vision system.

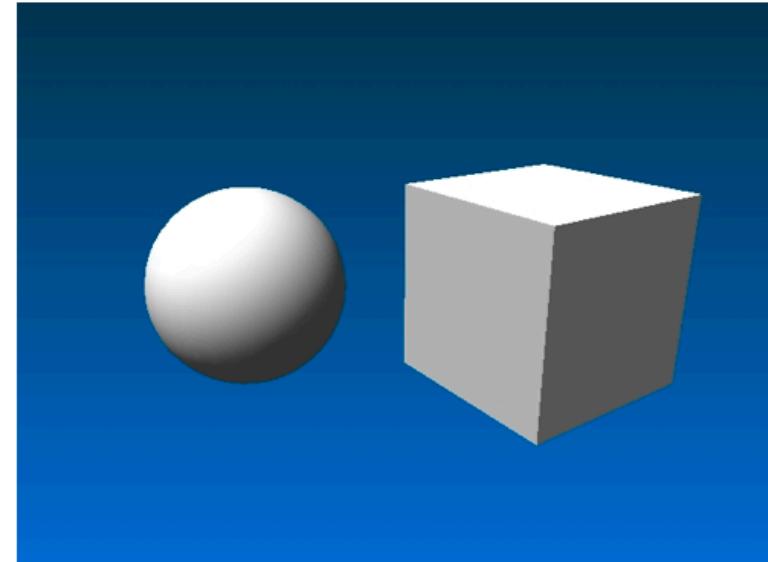
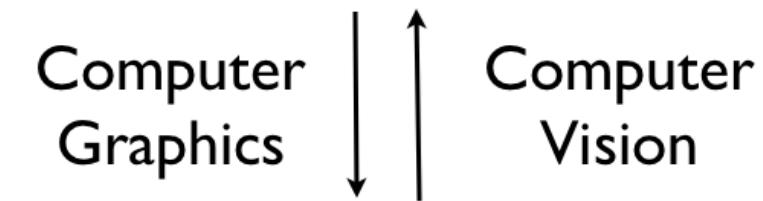
Interdisciplinary Nature of Computer Vision



Vision and Graphics

- **Graphics:**
 - go from the parameter space to the image space (rendering)
- **Vision:**
 - inverse graphics
 - more ill-posed
 - arguably harder

(cube, size, x_0 , y_0 , z_0 , θ_{xy} , θ_{xz} , θ_{yz} , ...)
(sphere, radius, x_1 , y_1 , z_1 , ...)



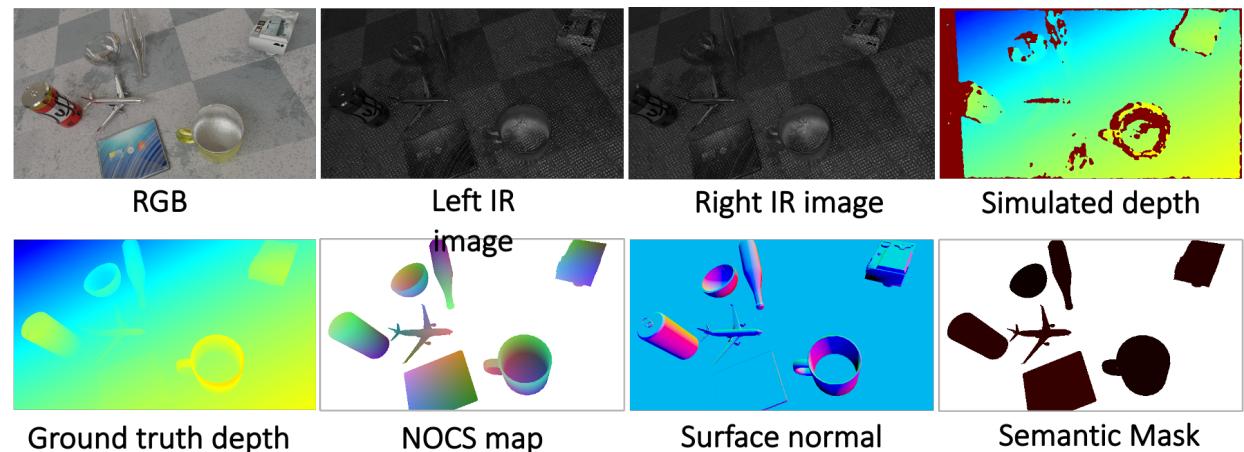
Leveraging Graphics for Vision

Synthetic data comes with **free** labels!

For outdoor semantic segmentation



For table-top depth/pose/surface normal prediction (ECCV 2022 from EPIC Lab)



Vision does Graphics Job

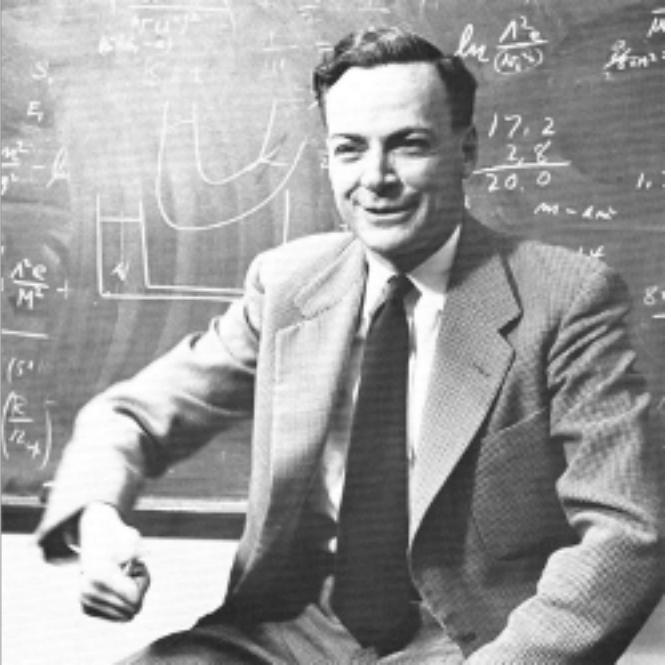
Training network to reproduce all input views of the scene



Neural Radiance Field (NeRF)

A brief intro to NeRF: <https://www.youtube.com/watch?v=JuH79E8rdKc>

Visual Content Generation



*What I cannot create,
I do not understand.*

- Generation or imagination is not a core function of human visual system.
- Richard Feynman: “What I cannot create, I do not understand”
- Thus, computer vision also deals with generation.

Applications: Human Face Generation



StyleGAN for facial image generation

Karras, Tero, Samuli Laine, and Timo Aila. "A style-based generator architecture for generative adversarial networks." *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2019.

Applications: Facial Reenactment

Animating Faces

A single model animates all images given only a single source image



Siarohin, A., Lathuilière, S., Tulyakov, S., Ricci, E., & Sebe, N. (2019). First order motion model for image animation. *Advances in Neural Information Processing Systems*, 32.

Applications: Style Transfer

Content target



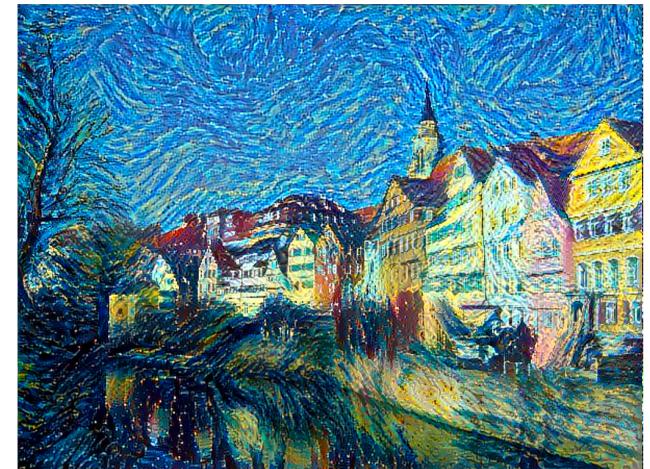
Style reference



+

=

Combination image



Neural Style Transfer

<https://gitee.com/happyjoejoe/deep-learning-with-python-notebooks/blob/master/8.3-neural-style-transfer.ipynb>

Applications: Text-to-Image Diffusion Model



A brain riding a rocketship heading towards the moon.



A dragon fruit wearing karate belt in the snow.



A small cactus wearing a straw hat and neon sunglasses in the Sahara desert.



A photo of a Corgi dog riding a bike in Times Square. It is wearing sunglasses and a beach hat.



A blue jay standing on a large basket of rainbow macarons.



The Toronto skyline with Google brain logo written in fireworks.



A bucket bag made of blue suede. The bag is decorated with intricate golden paisley patterns. The handle of the bag is made of rubies and pearls.



A single beam of light enter the room from the ceiling. The beam of light is illuminating an easel. On the easel there is a Rembrandt painting of a raccoon.

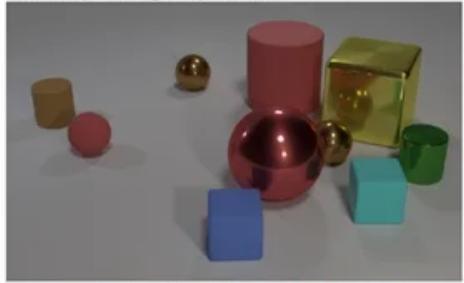
Imagen: <https://imagen.research.google/>

More Vision Language Tasks

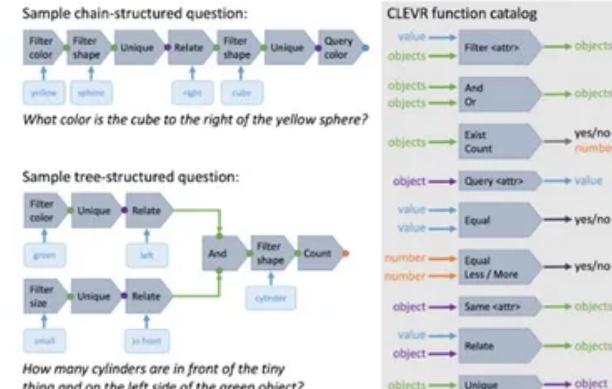
- Referring Expression
 - ReferIt Game, *EMNLP 2014*
 - RefCOCO, *ECCV 2016*
 - GuessWhat?!, *CVPR 2017*
- Visual Dialog
 - VisDial, *CVPR 2017*
 - Image Grounded Conversation, *ACL 2017*
 - Dialog-based Image Retrieval, *NIPS 2018*
- Text 2 image/video



Questions in CLEVR test various aspects of visual reasoning including **attribute identification**, **counting**, **comparison**, **spatial relationships**, and **logical operations**.



- Q: Are there an **equal number** of **large things** and **metal spheres**?
 Q: **What size** is the **cylinder** that is **left** of the **brown metal** thing that is **left** of the **big sphere**?
 Q: There is a **sphere** with the **same size** as the **metal cube**; is it **made of the same material** as the **small red sphere**?
 Q: **How many** objects are **either small cylinders** or **red things**?



Johnson, Justin, Bharath Hariharan, Laurens van der Maaten, Fei Fei Li, C. Lawrence Zitnick, and Ross Girshick. "Clevr: A diagnostic dataset for compositional language and elementary visual reasoning." In *CVPR*. 2017.

Intersection between computer vision and natural language processing (NLP)

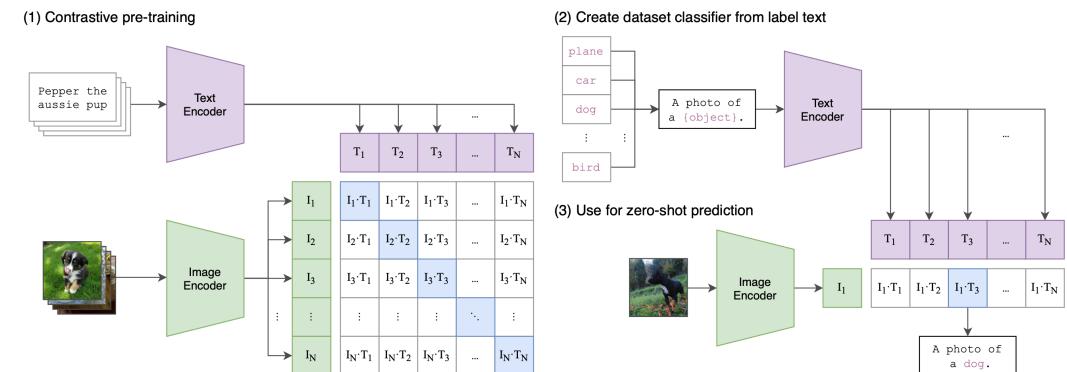


Figure 1. Summary of our approach. While standard image models jointly train an image feature extractor and a linear classifier to predict some label, CLIP jointly trains an image encoder and a text encoder to predict the correct pairings of a batch of (image, text) training examples. At test time the learned text encoder synthesizes a zero-shot linear classifier by embedding the names or descriptions of the target dataset's classes.

What Else in Computer Vision?

- Computer vision deals with
 - acquiring
 - processing and analyzing
 - understanding
 - generating or imagining
- visual data, and, what's more,
- providing visual feedbacks for body motions
 - helping making decisions
- for **embodied** agents.



Introduction to Computer Vision

Next week: Lecture 2,
Classic Vision I