# Robotics

Miao Li

Fall 2023, Wuhan University
WeChat: 15527576906
Email: limiao712@gmail.com

# Goal for this course

- **Design：soft hand design  x1**

- **Perception: vision, point cloud, tactile, force/torque x1**

- **Planning: sampling-based, optimization-based, learning-based x3**

- **Control: feedback, multi-modal x2**

- **Learning: imitation learning, RL x2**

- **Simulation tool (pybullet, matlab, OpenRAVE, Issac Nvidia, Gazebo)**

- **How to get a robot moving!**

# Today's Agenda

- **What is robot perception? (~12)**

- **Robot vision and computer vision (~5)**

- **Force sensing (~5)**

- **Tactile sensing (~5)**

- **Challenges of robot perception (10)**

- **Algorithms for perception**
  - **State estimation (~5)**
  - **End to end learning (~5)**
  - **Active perception (~5)**

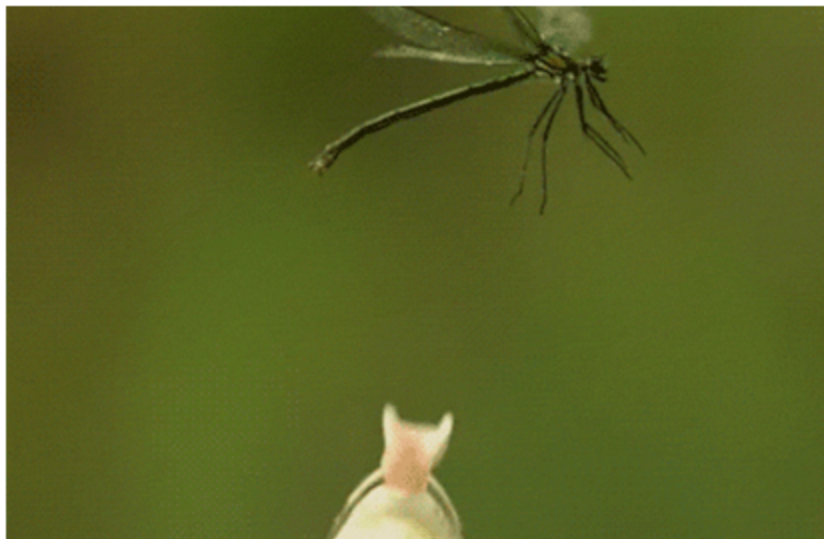- **Quick Review of Deep Learning (~20)**

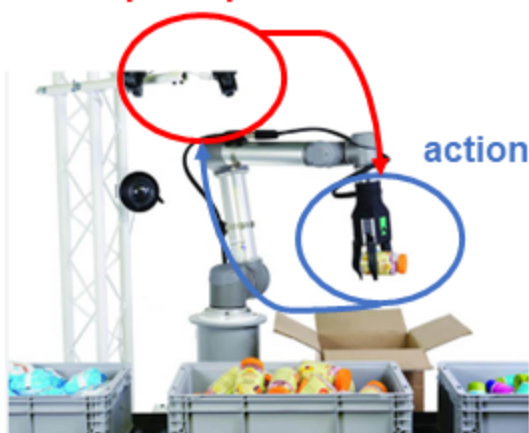# Incredible human skills

# Incredible animal skills

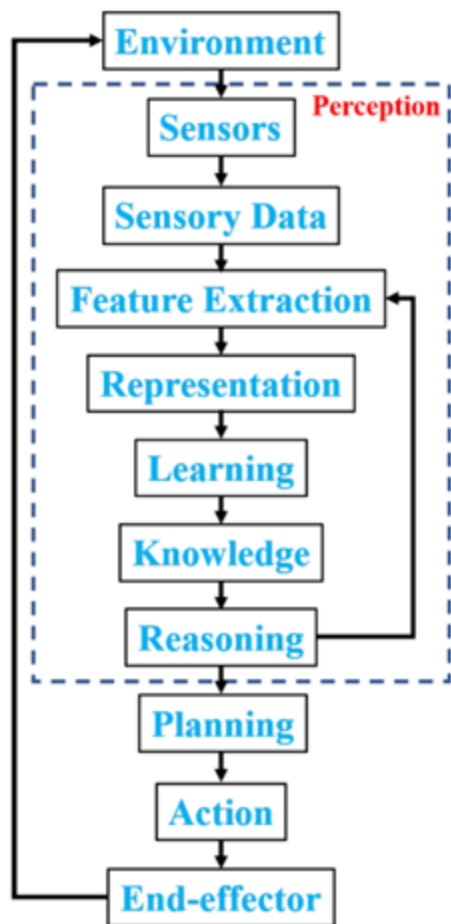**Robotics – Learn the mapping from perception to action**

Optimus is now capable of self-calibrating its arms and legs

Tesla's Optimus Robot Sort Objects Autonomously
https://www.youtube.com/watch?v=oL5YNtDUQXU&ab_channel=CNETHighlights

# Robotics – Learn the mapping from perception to action

Robotics – Learn the mapping from **perception** to action

# Why robot perception?

**Making sense of the unstructured environment …**
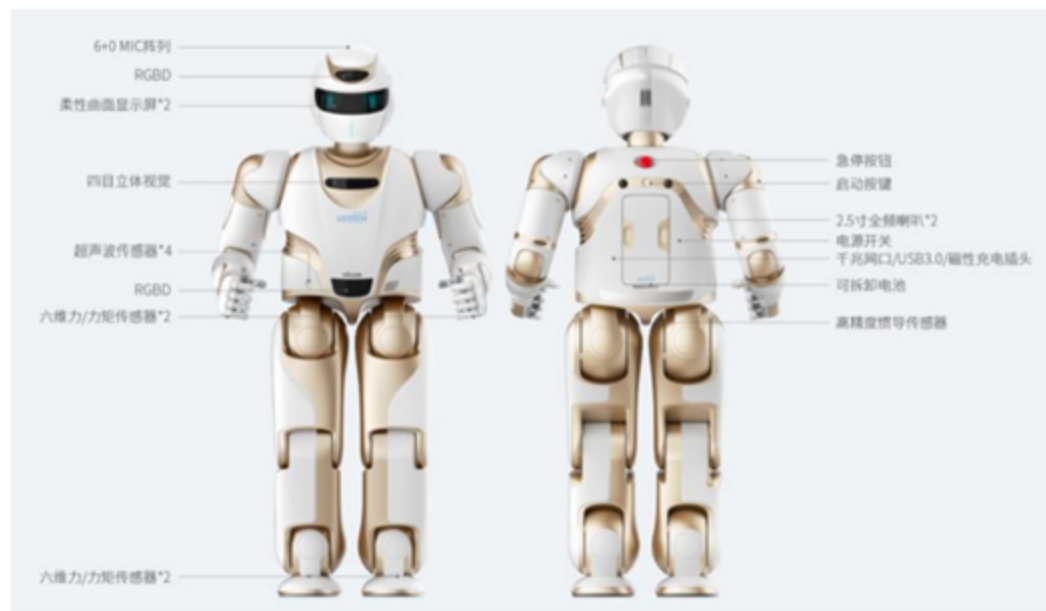


- Incomplete knowledge of the scene

- Imperfect actions may lead to failure
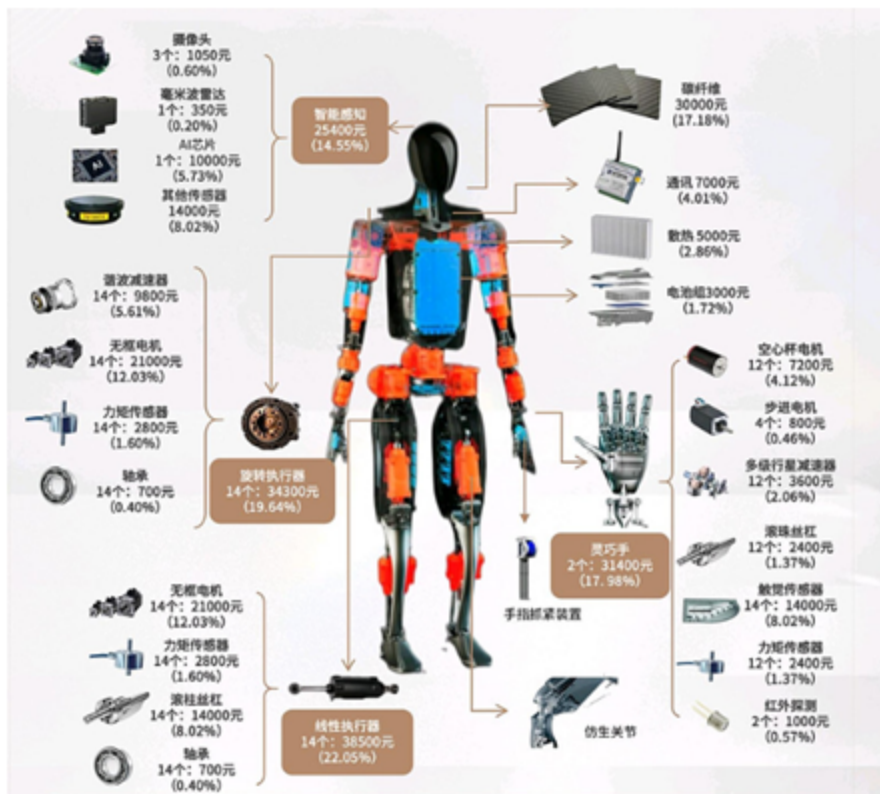
- Environment dynamics

# Robot sensor

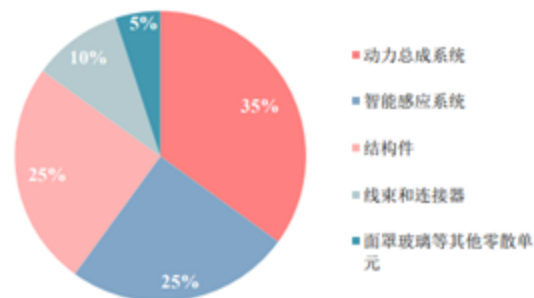## Understand the real world through different sensors

# Robot sensor

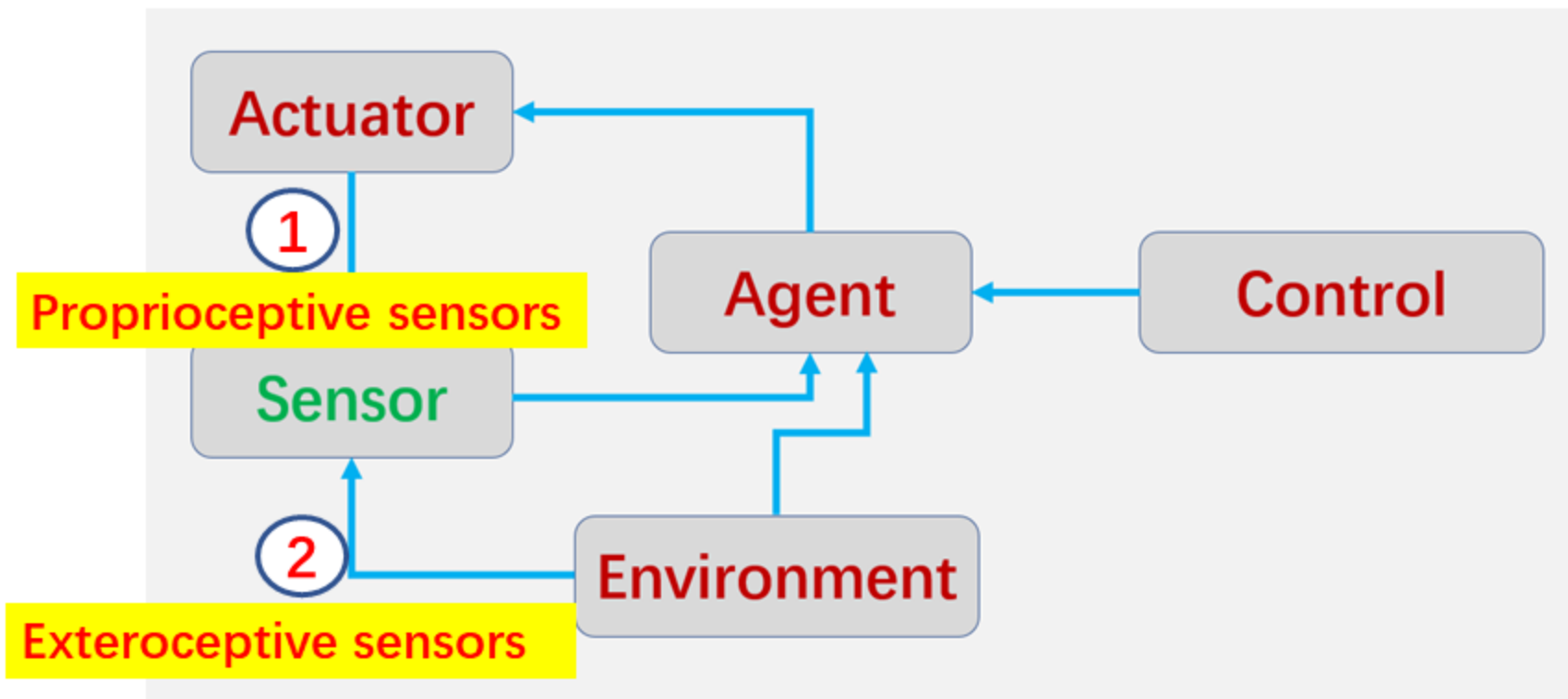## Understand the real world through different sensors

# Robot sensor

# Robot sensor

- **proprioceptive sensors** measure the internal state of the robot (position and velocity of joints, but also torque at joints or acceleration of links)
  - kinematic calibration, identification of dynamic parameters, control
- **exteroceptive sensors** measure/characterize robot interaction with the environment, enhancing its autonomy(forces/torques, proximity, **vision**, but also sensors for sound, smoke humidity, ...)
  - control of interaction with the environment, obstacle avoidance localization of mobile robots, navigation in unknown environments

# Robot vision vs. Computer vision

**Computer vision** tasks include methods for <u>acquiring</u>, <u>processing</u>, <u>analyzing</u> and understanding digital images, and extraction of <u>high-dimensional</u> data from the real world in order to produce numerical or symbolic information, e.g. in the forms of decision.
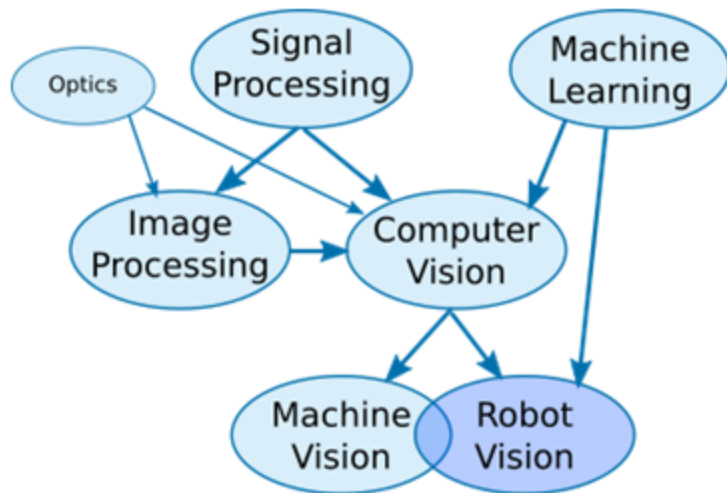




**Deep Learning**

# Robot vision vs. Computer vision

| Technique | Input | Output |
|---|---|---|
| Signal Processing | Electrical signals | Electrical signals |
| Image Processing | Images | Images |
| Computer Vision | Images | Information/features |
| Pattern Recognition/Machine Learning | Information/features | Information |
| Machine Vision | Images | Information |
| Robot Vision | Images | Physical Action |



https://blog.robotiq.com/robot-vision-vs-computer-vision-whats-the-difference
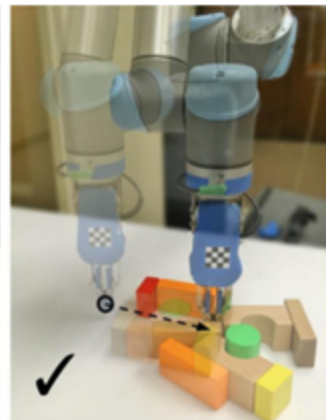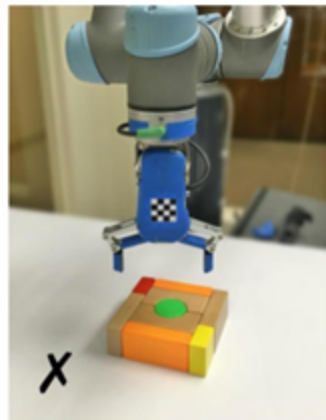
# Robot vision vs. Computer vision

- Robot vision is **embodied**, **active**, and environmentally **situated**.

- **Embodied**: Robots have physical bodies and experience the world directly. Their actions are part of a dynamic with the world and have immediate feedback on their own sensation.

- **Active**: Robots are active perceivers. It knows why it wishes to sense, and chooses what to perceive, and determines how, when and where to achieve that perception.

- **Situated**: Robots are situated in the world. They do not deal with abstract descriptions, but with the here and now of the world directly influencing the behaviorof the system.

# Robot vision vs. Computer vision



[Levine et al., IJRR 2016]

[Zeng et al., IROS 2018]

# 2D camera

1963 – Lawrence Roberts, the Father of Computer Vision publishes "Machine Perception Of Three-Dimensional Solids" where he discusses extracting 3D information about solid objects from 2D images. This lead to much research in MIT's artificial intelligence lab and other research institutions looking at computer vision in the context of blocks and simple objects.

1966 – The summer project at MIT marks the landmark in the development of pattern recognition.



MASSACHUSETTS INSTITUTE OF TECHNOLOGY
PROJECT MAC

Artificial Intelligence Group                July 7, 1966
Vision Memo. No. 100.


THE SUMMER VISION PROJECT

Seymour Papert

https://emergentvisiontec.com/tech-portal/evolution-of-machine-vision/

# 2D camera

The Optical Path from the Lens Through the Mirror to the Viewfinder

Pentaprism or pentamirror

Shutter-release button

Viewfinder screen

Viewfinder eyepiece window

Image sensor

Shutter

Lens

Aperture

Mirror

# 2D camera



different ways to scan an **array of pixels**

**CAMERA**

SYNCHRONIZATION

ANALOG ELECTRONICS

*analog video signal*

LIGHT

shutter

sensor

lens

FG DSP CPU

Frame Grabber (A/D conversion)

dedicated boards (low-level vision)

high-level vision

VIDEO STANDARDS
CCIR (Europe, Australia): 625 lines, 25 Hz
RS170 (USA, Japan): 525 lines, 30 Hz
video signal = 1 V peak-to-peak

*it becomes then a vision-driven robot control systems ... (with real-time constraints!)*

robot

controller

# 2D camera

**Visual servoing**, also known as **vision-based robot control** and abbreviated **VS**, is a technique which uses feedback information extracted from a vision sensor (visual feedback) to control the motion of a robot. One of the earliest papers that talks about visual servoing was from the SRI International Labs in 1979.



1. *"Basic Concept and Technical Terms"*. *Ishikawa Watanabe Laboratory, University of Tokyo. Retrieved 12 February 2015.*
2.^ Agin, G.J., "Real Time Control of a Robot with a Mobile Camera". Technical Note 179, SRI International, Feb. 1979.

# 2D image format



What the computer sees

An image is just a matrix of numbers [0,255]!
i.e., 1080x1080x3 for an RGB image

Images Are Numbers

# 3D from stereo



Object

Point in 3D
(on object)

Point projection    Baseline    Point projection

left image                    right image

disparity: the difference in image location of the same 3D
point when projected under perspective to two different cameras.

$$d = xleft - xright$$

# 3D from stereo

image plane

z

Z

camera L

f

xl

baseline → b

f

camera R

xr

x-b

P=(x,z)

X

$$\frac{z}{f} = \frac{x}{xl}$$

$$\frac{z}{f} = \frac{x-b}{xr}$$

$$\frac{z}{f} = \frac{y}{yl} = \frac{y}{yr}$$

y-axis is
perpendicular
to the page.

# 3D from stereo

For stereo cameras with parallel optical axes, focal length f, baseline b, corresponding image points (xl,yl) and (xr,yr) with disparity d:

$$z = f*b / (xl - xr) = f*b/d$$

$$x = xl*z/f \quad or \quad b + xr*z/f$$

$$y = yl*z/f \quad or \quad yr*z/f$$

This method of determining depth from disparity is called **triangulation.**

# 3D from stereo



left image        right image

We need to find the "correspondence".

# 3D from structure

3D data can also be derived using

- a single camera

- a light source that can produce
  stripe(s) on the 3D object

light stripe

light
source

camera

# 3D from structure

3D data can also be derived using

- a single camera

- a light source that can produce
  stripe(s) on the 3D object

$$[x \ y \ z] = \frac{b}{f \cot \theta - x'} [x' \ y' \ f]$$

3D        image



3D point
(x, y, z)

θ

(0,0,0)

light
source

b    f

x axis

(x',y',f)

$$\frac{x}{x'} = \frac{(b+x) \cdot \tan\theta}{f}$$

# 3D vision format



- PLY - a polygon file format, developed at Stanford University by Turk et al

- STL - a file format native to the stereolithography CAD software created by 3D Systems

- OBJ - a geometry definition file format first developed by Wavefront Technologies

- X3D - the ISO standard XML-based file format for representing 3D computer graphics data

- and many others

https://pointclouds.org/documentation/tutorials/pcd_file_format.html

# 3D camera application



Source: Keyence website

# 3D camera application

# F/T sensor principle

- indirect information obtained from the measure of <span style="color:red">deformation</span> of an elastic element subject to the force or torque to be measured

- basic component is a *strain gauge*: uses the variation of the resistance R of a metal conductor when its length L or cross-section S vary
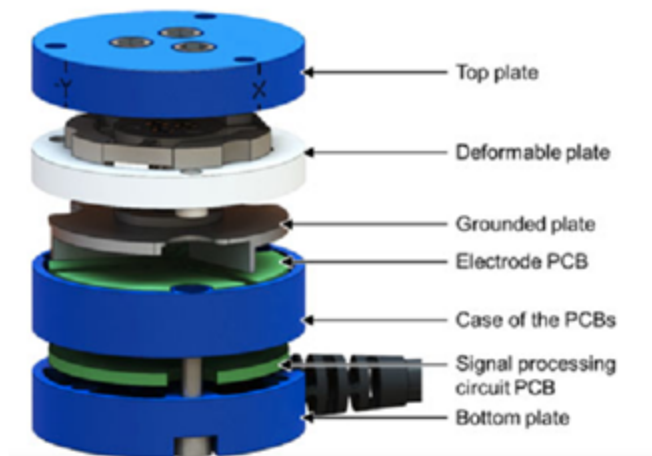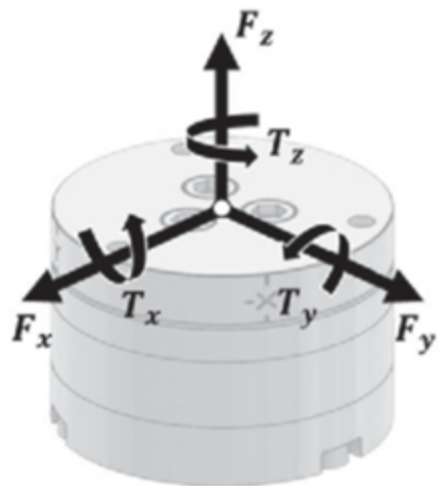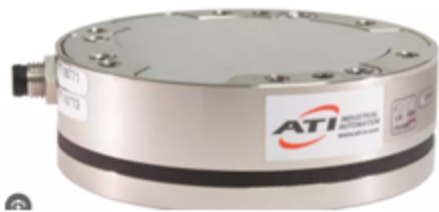
$$\frac{\partial R}{\partial L} > 0 \qquad \frac{\partial R}{\partial S} < 0$$

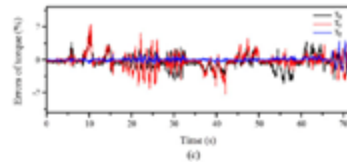$$\frac{\partial R}{\partial T} \text{ small}$$

temperature



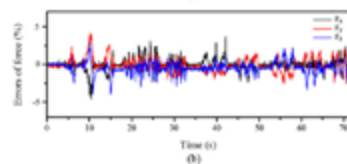Strain sensitive pattern

Terminals

Tension: area narrows, resistance increases.

Higher resistance

Compression: area thickens, resistance decreases.
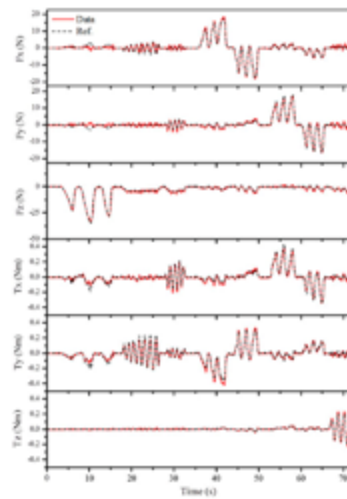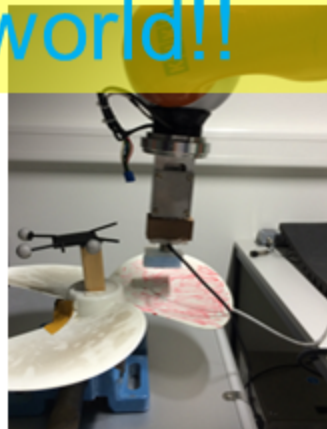
Lower resistance

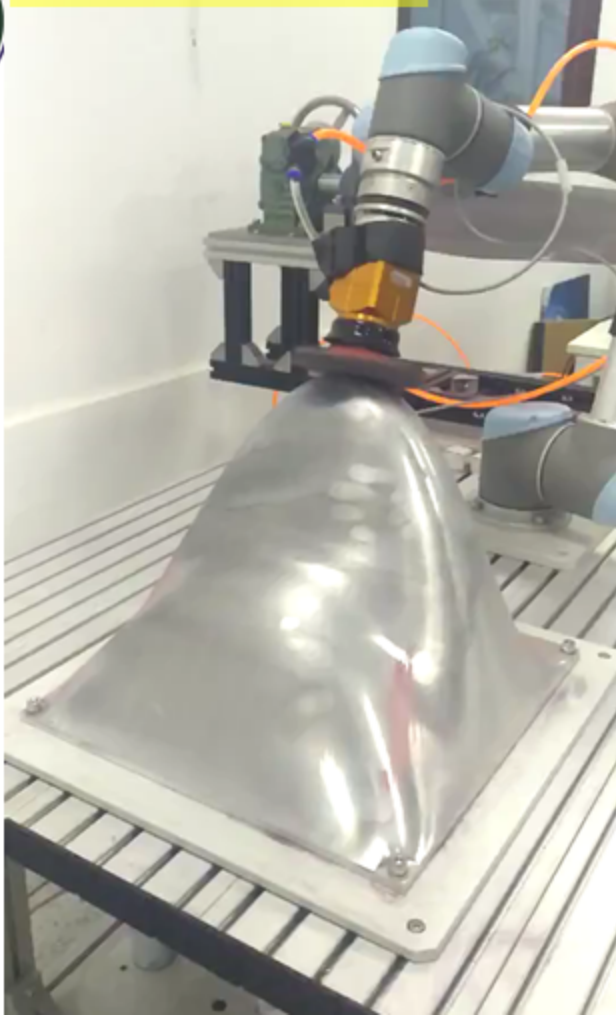# F/T sensor principle

# F/T information format

# F/T sensor application



Force is an essential information for the robots to physically interact with the world!!
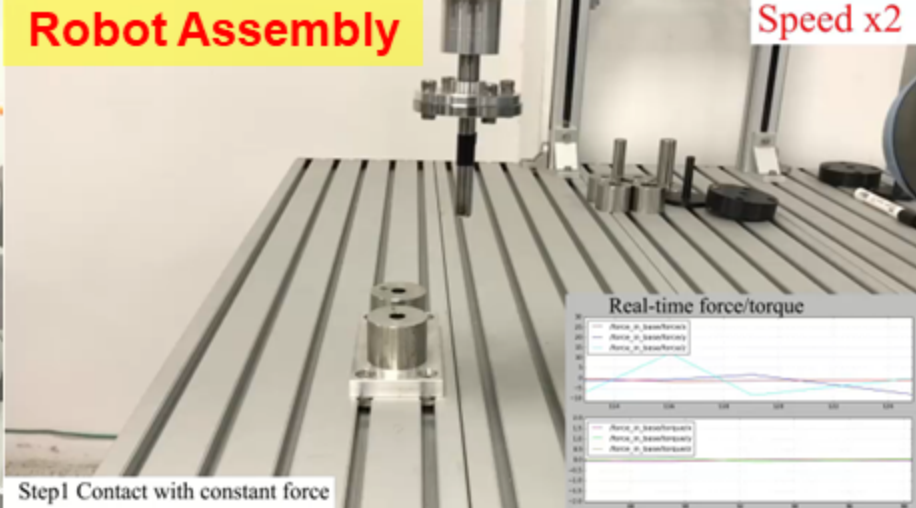
# Tactile sensor principle

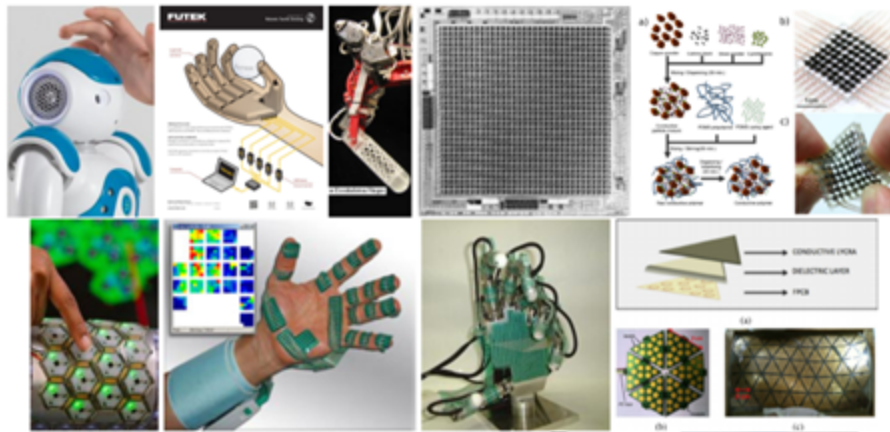# Tactile sensor history



**Large-Area Soft e-Skin: The Challenges Beyond Sensor Designs**

# Tactile sensor applications

# Tactile sensor principle



**Fluid-type touchpad**

Colored water (translucent red)

Transparent acrylic plate

Two-layered elastic membrane (black and white)

Metallic plate

CCD camera

LED lights

One-way mirror

Touchpad

Object

Elastic cover

Optical waveguide

Scattered light

Optical detector

Processor

Light source

Object

Optical fiber

Sensing surface

Reflective layer

Sensing element

Light emitter

# Robot Skin

# Why is "optical skin" a good idea?

- Easy to get millions of sensors (pixels)

- Bonus: Proximity Sense.

- Separate deformation and its measurement.

- Outer layer has nothing in it. Cheap to repair or replace. Can be optimized for desired mechanical and lifetime properties.

- Connections (solder joints etc.) and wires aren't deforming.

- Reliability – fewer wires, components. Cameras are already reliable.

- Whole-body vision: reduce occlusion.

- Avoid rigid printed circuit boards or chips.

# Finger Vision: Proximity Sensing Seeing with your fingers

# Whole Body Vision: Origins

Argus (Greece)

Hyakume (Japan)

# Goal and Challenge

- Can we build 100 camera system on full robot
- with full multimodal sensor suite
- and networking (GigE?)
- and processing (NVIDIA Tegra?)
- and reasonable power budget
- and do something interesting (Repair? Cooking?)

Blue: cameras
Red: structured light sources
or DLP projectors

# Some reference books on sensors using in mechatronics, robotics, vision control···

# Today's Agenda

- **What is robot perception? (~12)**

- **Robot vision and computer vision (~5)**

- **Force sensing (~5)**

- **Tactile sensing (~5)**

- **Challenges of robot perception (10)**

- **Algorithms for perception**

  - **State estimation  (~5)**

  - **End to end learning  (~5)**

  - **Active perception  (~5)**

- **Quick Review of Deep Learning (~20)**

# Why robot perception is difficult?

1. **Uncertainty**: noise is everywhere!

2. **Modalities**: neural network architectures designed for different sensory modalities

3. **Representations**: representation learning algorithms without strong supervision

4. **Tasks**: state estimation tasks for robot navigation and manipulation

5. **Embodiment**: active perception for embodied visual intelligence

# Why do we need a filter?



**Noise is everywhere!**

- - - original coordinate
——— sensor's readings
——— filtered by Kalman

- The main reason to filter a signal is to reduce and smooth out **high-frequency noise** associated with a measurement such as flow, pressure, level or temperature
- When a noisy signal is used in control, filtering is important for effective derivative action and for avoiding excessive movement in the controller output that causes valve wear or disturbs other control loops.
- Ideally, we want to estimate the underlying signal without noise, introducing as little distortion as possible.

https://en.wikipedia.org/wiki/Smoothing



Engineers use filtering to extract the useful information from noisy signals.

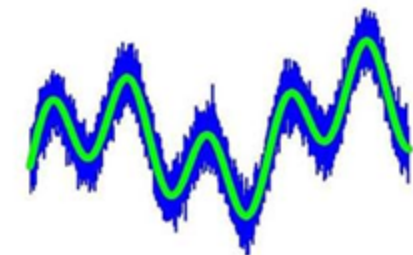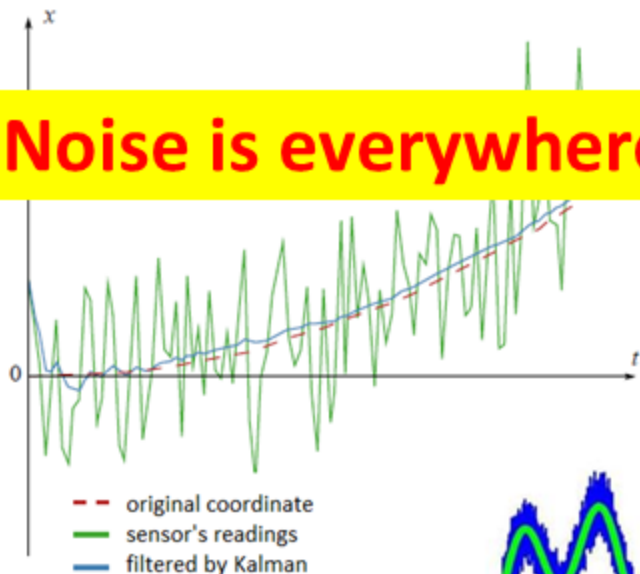| Algorithm | Overview and uses | Pros |
|---|---|---|
| Additive smoothing | used to smooth categorical data. | |
| Butterworth filter | Slower roll-off than a Chebyshev Type I/Type II filter or an elliptic filter | • More linear phase response in the pass-band than Chebyshev Type I/Type II and elliptic filters can achieve.<br>• Designed to have a frequency response as flat as possible in the passband. |
| Chebyshev filter | Has a steeper roll-off and more passband ripple (type I) or stopband ripple (type II) than Butterworth filters. | • Minimizes the error between the idealized and the actual filter characteristic over the range of the filter |
| Digital filter | Used on a sampled, discrete-time signal to reduce or enhance certain aspects of that signal | |
| Elliptic filter | | |
| Exponential smoothing | • Used to reduce irregularities (random fluctuations) in time series data, thus providing a clearer view of the true underlying behaviour of the series.<br>• Also, provides an effective means of predicting future values of the time series (forecasting). | |
| Kalman filter | • Uses a series of measurements observed over time, containing statistical noise and other inaccuracies by estimating a joint probability distribution over the variables for each timeframe. | Estimates of unknown variables it produces tend to be more accurate than those based on a single measurement alone |
| Kernel smoother | • used to estimate a real valued function as the weighted average of neighboring observed data.<br>• most appropriate when the dimension of the predictor is low (p < 3), for example for data visualization. | The estimated function is smooth, and the level of smoothness is set by a single parameter. |

# Kalman Filter

**motion and sensing discrete-time model for estimation**

$$\boldsymbol{\xi}(k) = \begin{pmatrix} 1 & T \\ 0 & 1 \end{pmatrix} \boldsymbol{\xi}(k-1) + \boldsymbol{\mu}$$

$$z(k) = \begin{pmatrix} 1 & 0 \end{pmatrix} \boldsymbol{\xi}(k) + \nu$$

zero mean Gaussian noises with (co)variances $Q$ (a matrix) and $R$

noisy **position measure** (encoder output)

$T$ = sampling time

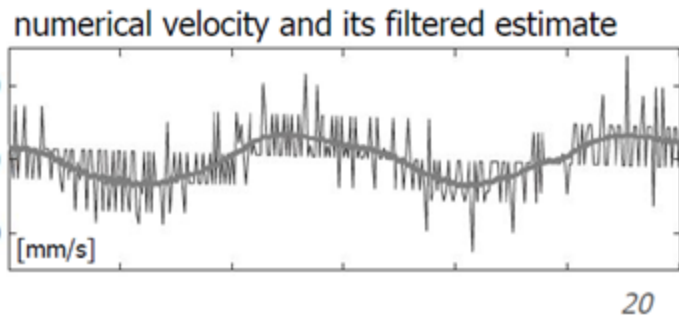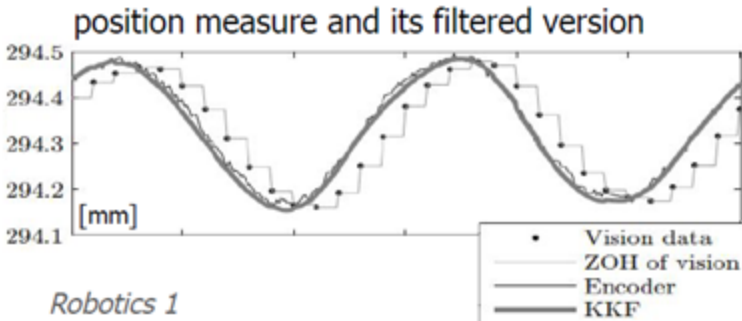$$\boldsymbol{\xi}(k) = (x(k)\ \dot{x}(k))^T$$

actual state

unmeasured velocity

design a (linear) **Kalman filter** providing an **estimate** $\hat{\boldsymbol{\xi}}(k)$ of the model state

$$\hat{\boldsymbol{\xi}}(k) = \begin{pmatrix} 1 & T \\ 0 & 1 \end{pmatrix} \hat{\boldsymbol{\xi}}(k-1) + \boldsymbol{K}_k \left( z(k) - \begin{pmatrix} 1 & 0 \end{pmatrix} \begin{pmatrix} 1 & T \\ 0 & 1 \end{pmatrix} \hat{\boldsymbol{\xi}}(k-1) \right)$$

using the **optimal** Kalman gain $\boldsymbol{K}_k$

(a priori) **prediction**       **correction** (based on the measured output)

position measure and its filtered version

numerical velocity and its filtered estimate



- Vision data
- ZOH of vision
- Encoder
- KKF

# Robot Perception: Modality

Pixels (from RGB cameras)

What the computer sees

An image is just a big grid of numbers between [0, 255]:

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

$(x_1, y_1, z_1)$

$(x_2, y_2, z_2)$

[Source: PointNet++; Qi et al. 2016]

Point cloud (from structure sensors)

Reaching    Alignment    Insertion

$F_z (N)$

Time (ms)

[Source: Lee*, Zhu*, et al. 2018]

Time series (from F/T sensors)

$a_2$

$a_1$

[Source: Calandra et al. 2018]

Tactile data (from the GelSights sensors)

# Robot Perception: Modality



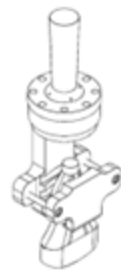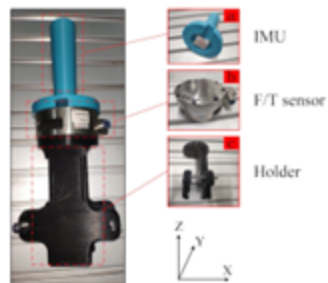Sonographer ultrasound scanning process     CAD model of probe holder     Probe holder with sensors
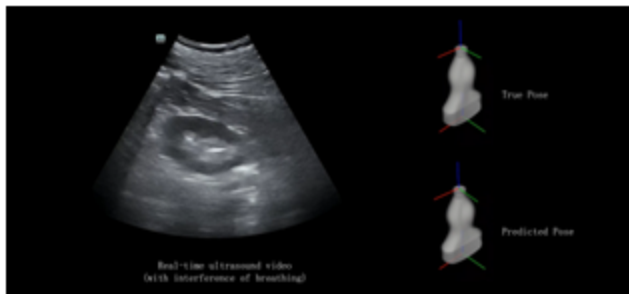
# Robot Perception: Modality

# Robot Perception: Modality



Electroencephalogram (EEG)

Electromyography (EMG)

Gaze tracking

Environment

Heart rate

Motion detection

Blood Pressure / Temperature

Sleep patterns

## Different Types of Sensors

Thermistor (Temperature Sensor)

IR Sensor (Transmissive Type)

IR Sensor (Reflective Type)

Ultrasonic Sensor

Gyroscope Sensor

Accelerometer Sensor

Rain Sensor

Soil Moisture Sensor

Phototransistor (Light Sensor)

Water Flow Sensor

Heartbeat Sensor

Alcohol Sensor

Color Sensor

PIR Sensor

Gas Sensor

Smoke Sensor

LM35 (Temperature Sensor)

IR Receiver

LDR (Light Sensor)

www.electricaltechnology.org

Humidity Sensor

Flex Sensor

Touch Sensor

Solar Cell Light Sensor

Metal Dedector

Real Time Clock Sensor

Vibration Sensor

# Robot Perception: Modality



```
Environment
    ↓
  Sensors          Perception
    ↓
Sensory Data
    ↓
Feature Extraction
    ↓
Representation
    ↓
  Learning
    ↓
 Knowledge
    ↓
  Reasoning
    ↓
  Planning
    ↓
   Action
    ↓
 End-effector
```
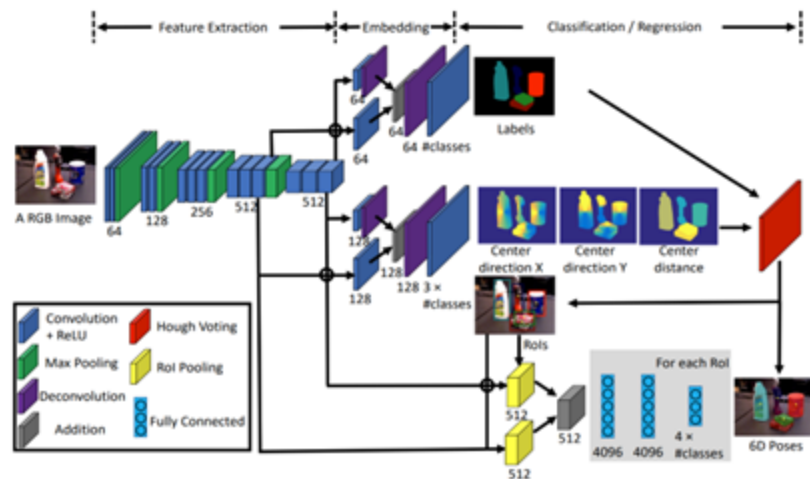
●**How can we design the algorithms (neural networks) that can effectively process the raw sensory data in different forms?**

# Robot Perception: Modality



PoseCNN: A Convolutional Neural Network for 6D
Object Pose Estimation in Cluttered Scenes

Yu Xiang[1,2], Tanner Schmidt[2], Venkatraman Narayanan[3] and Dieter Fox[1,2]
[1]NVIDIA Research, [2]University of Washington, [3]Carnegie Mellon University
yux@nvidia.com, tws10@cs.washington.edu, venkatraman@cs.cmu.edu, dieterf@nvidia.com

# Robot Perception: Modality



Environment

Sensors — Perception

Sensory Data

Feature Extraction

Representation

Learning

Knowledge

Reasoning

Planning

Action

End-effector

PointNet

mug?
table?
car?

Classification    Part Segmentation    Semantic Segmentation

*Classification Network*
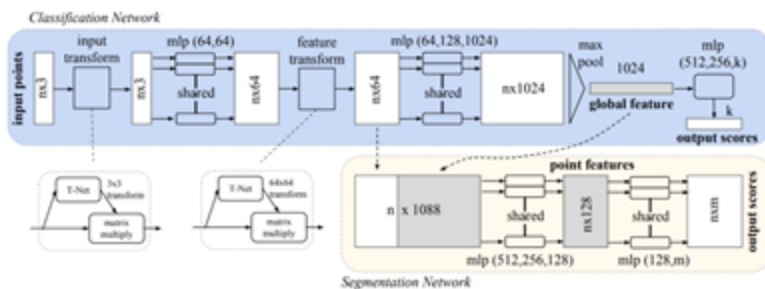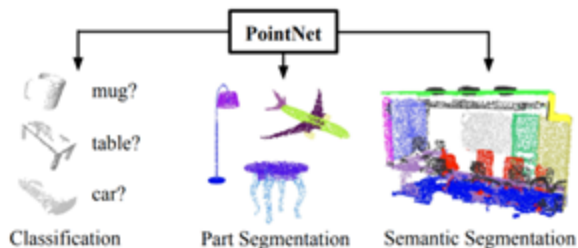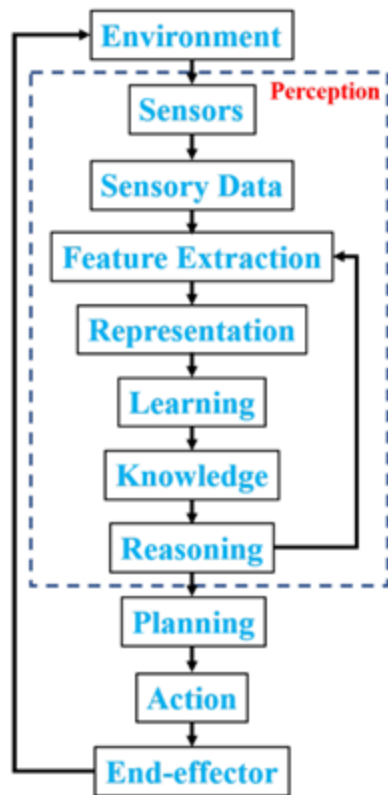
*Segmentation Network*

Figure 2. **PointNet Architecture.** The classification network takes *n* points as input, applies input and feature transformations, and then aggregates point features by max pooling. The output is classification scores for *k* classes. The segmentation network is an extension to the classification net. It concatenates global and local features and outputs per point scores. "mlp" stands for multi-layer perceptron, numbers in bracket are layer sizes. Batchnorm is used for all layers with ReLU. Dropout layers are used for the last mlp in classification net.

**PointNet: Deep Learning on Point Sets for 3D Classification and Segmentation**

Charles R. Qi*    Hao Su*    Kaichun Mo    Leonidas J. Guibas
Stanford University

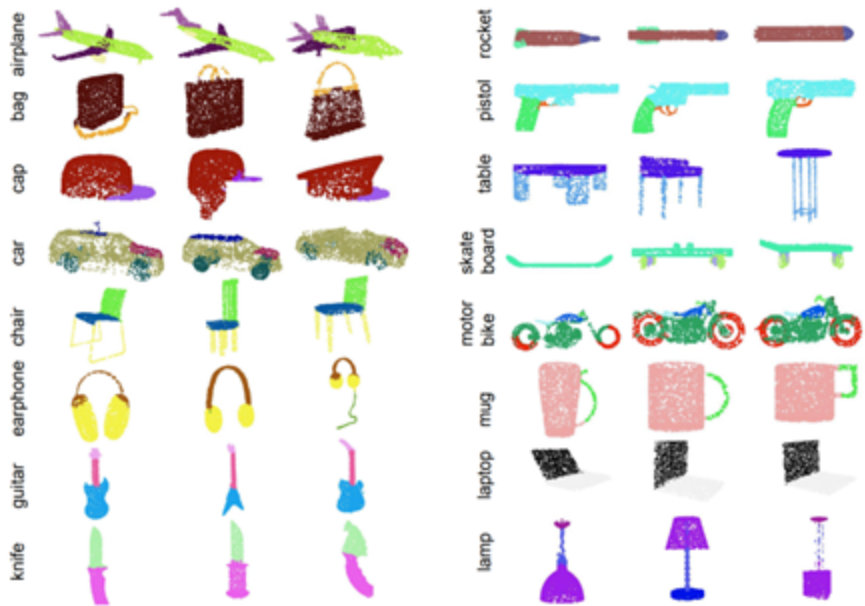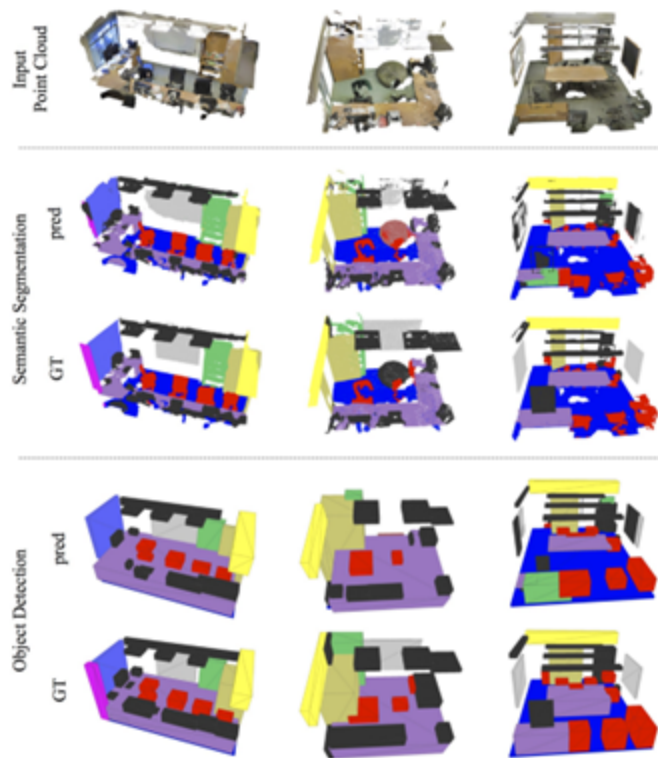# Robot Perception: Modality



Figure 21. **PointNet segmentation results on complete CAD models.**

# Robot Perception: Representation

A fundamental problem in robot perception is to learn the proper **representations** of the unstructured world.



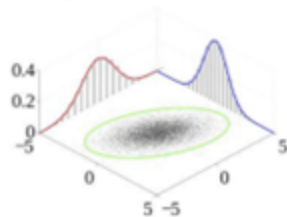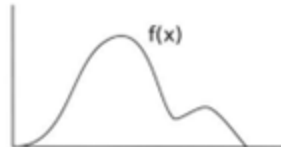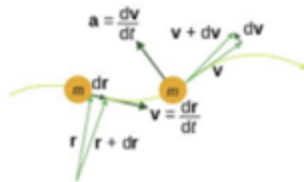[Source: Stanford CS331b]

# Robot Perception: Representation

"Solving a problem simply means representing it so as to make the solution transparent."

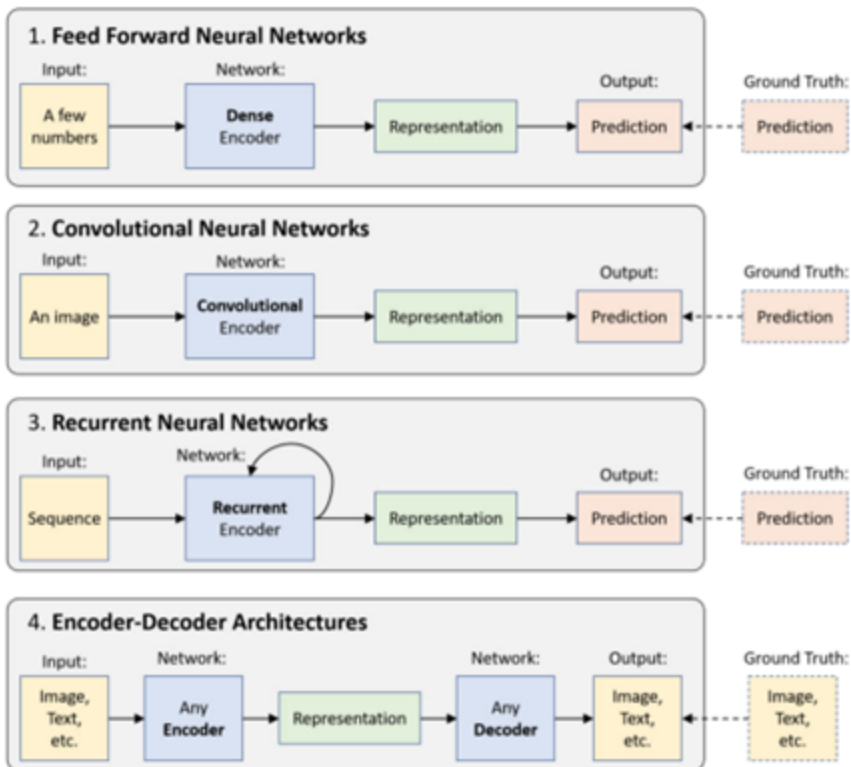Herbert A. Simon, Sciences of the Artificial

# Robot Perception: Representation



**Supervised Learning**

**1. Feed Forward Neural Networks**

Input: A few numbers → Network: Dense Encoder → Representation → Output: Prediction ⇠ Ground Truth: Prediction

**2. Convolutional Neural Networks**

Input: An image → Network: Convolutional Encoder → Representation → Output: Prediction ⇠ Ground Truth: Prediction

**3. Recurrent Neural Networks**

Input: Sequence → Network: Recurrent Encoder → Representation → Output: Prediction ⇠ Ground Truth: Prediction

**4. Encoder-Decoder Architectures**

Input: Image, Text, etc. → Network: Any Encoder → Representation → Network: Any Decoder → Output: Image, Text, etc. ⇠ Ground Truth: Image, Text, etc.

**Unsupervised Learning**

**5. Autoencoder**

Input: Image, Text, etc. → Network: Any Encoder → Representation → Network: Any Decoder → Ground Truth: Exact copy of input

**6. Generative Adversarial Networks**

*Throw away after training*

Input: Noise → Network: Generator → Output: Fake Image → Network: Discriminator → Prediction: Real or Fake
Real Image

**Reinforcement Learning**

**7. Networks for Actions, Values, Policies, and Models**

Input: World State Sample → Network: Any Encoder → Representation → Output: Action ⇠ Ground Truth: Reward

[6.S094, MIT]

# Robot Perception: Representation

How can we learn representations that fuse **multiple sensory modalities** together?

**Next course will cover this part**



combining vision and force for manipulation

[Lee*, Zhu*, et al. 2018]

# Robot Perception: Representation



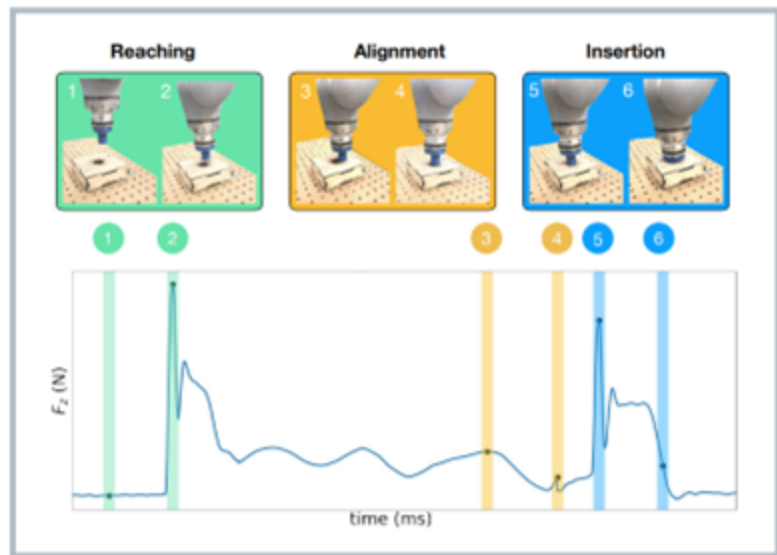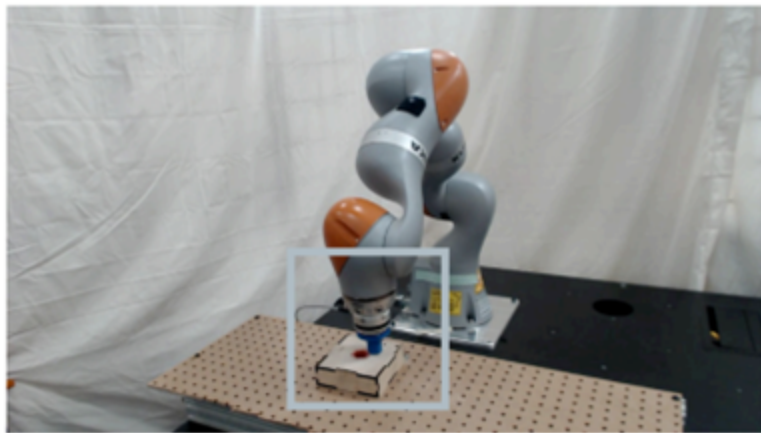(a) Three Poses Estimate
(b) SelectLSTM

PoseFusion: Robust Object-in-Hand Pose Estimation with SelectLSTM

Yuyang Tu[†1], Junnan Jiang[†2], Shuang Li[1], Norman Hendrich[1], Miao Li[*2] and Jianwei Zhang[*1]

PoseFusion: Robust Object-in-Hand Pose Estimation with SelectRNN

Yuyang Tu[†1], Junnan Jiang[†2], Shuang Li[1], Norman Hendrich[1], Miao Li[2] and Jianwei Zhang[1]

[1]Universität Hamburg, [2]Wuhan University, †denote equal contribution

与汉堡大学张建伟教授合作 (IROS 2023)

# Robot Perception: Representation



Robot Cooking with Stir-fry:
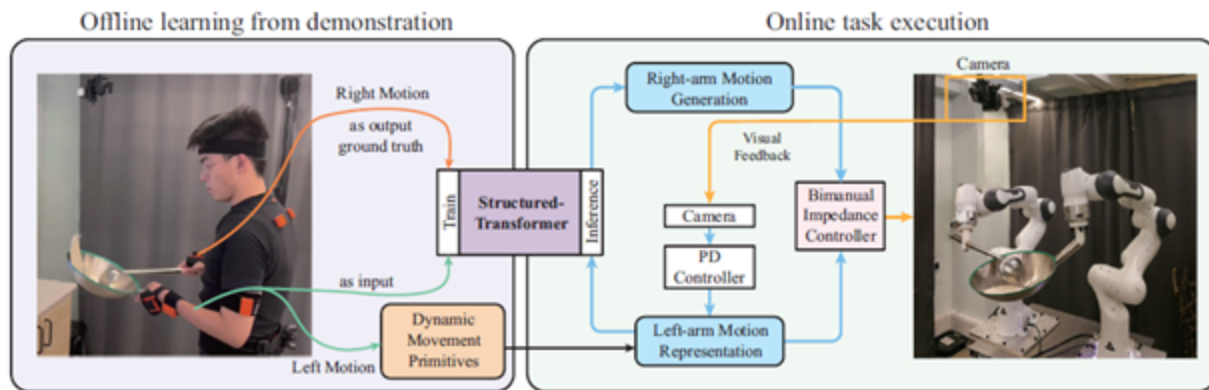Bimanual Non-prehensile Manipulation of Semi-fluid Objects

Junjia Liu[1], Yiting Chen[1,2], Zhipeng Dong[1], Shixiong Wang[1],
Sylvain Calinon[3], Miao Li[†4], and Fei Chen[†1], Senior Member, IEEE

与香港中文大学陈翡教授合作（RAL 2022）

# Robot Perception: Representation

# Robot Perception: Task



**Environment**

**Sensors** — Perception

**Sensory Data**

**Feature Extraction**

**Representation**

**Learning**

**Knowledge**

**Reasoning**

**Planning**

**Action**

**End-effector**

*Task ?*

*Task Evaluation.*

Noisy Sensory Data

State Representation

Perception & Computer Vision

Robot Control & Decision Making

**Robotic grasping, manipulation, navigation, exploration, …**

# Robot Perception: Task



Instruction:
I need to hammer a nail, what object from the scene might be useful?

Prediction:
Rocks. Action: 1 129 138 122 132 132 106 127

https://www.deepmind.com/blog/rt-2-new-model-translates-vision-and-language-into-action

**It is in general very difficult to represent a task, human use everyday language to express it.**

# Robot Perception: Embodiment

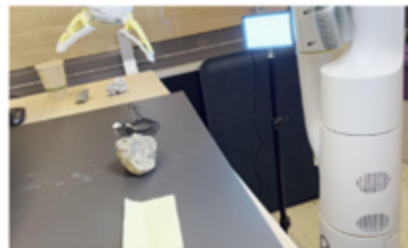Input-Output Picture (Susan Hurley, 1998)

**Conventional View of Perception**

- Perception is the process of building an internal representation of the environment

- Perception is input from world to mind, and action is output from mind to world, thought is the mediating process.

[Action in Perception, Alva Noë 2004]

"We see in order to move; we move in order to see." – William Gibson

# Robot Perception: Embodiment



Pebbles (James J. Gibson 1966)

**Embodied View of Perception**

- Subjects asked to find a reference object among a set of irregularly-shaped objects

- Three groups

  a. Passive observers of one static image (49%)

  b. Observers of moving shapes (72%)

  c. Interactive observers (99%)

- The ability to condition input signals with actions is crucial to perception.

Gibson, J. J. (1950). *The Perception of the Visual World*. Oxford England: Houghton Mifflin. ISBN 978-1114828087.

*Gibson, J. J. (1966). The senses considered as perceptual systems. Oxford England: Houghton Mifflin.* ISBN 978-0313239618.

# Robot Perception: Embodiment



View
Selection

scene observation completion

Where to
look next?

object observation completion

How to
manipulate?

[Ramakrishnan et al. 2019]

Physical
Interaction

Grasping   Pushing   Poking

Physical Interaction Data

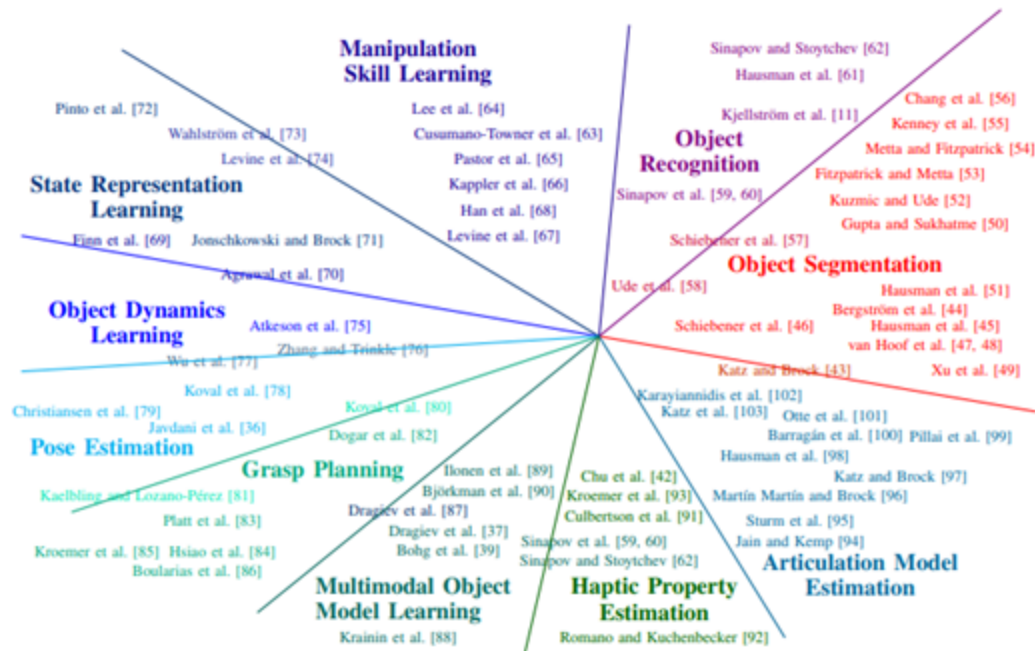Learned Visual Representation

[Pinto et al. 2016]

active perception
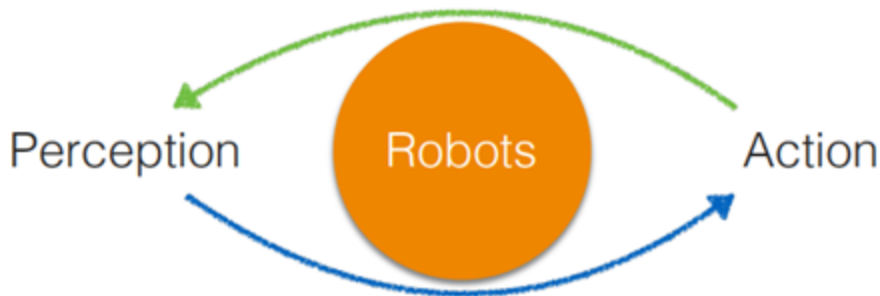
# Robot Perception: Embodiment



Interactive Perception: Leveraging Action in
Perception and Perception in Action

Jeannette Bohg*, *Member, IEEE*, Karol Hausman*, *Student Member, IEEE*, Bharath Sankaran*, *Student Member, IEEE*, Oliver Brock, *Senior Member, IEEE*, Danica Kragic, *Fellow, IEEE*, Stefan Schaal, *Fellow, IEEE*, and Gaurav Sukhatme, *Fellow, IEEE*

# Robot Perception: Embodiment



Perception — Robots — Action

How robots develop better perception from embodied sensorimotor experiences
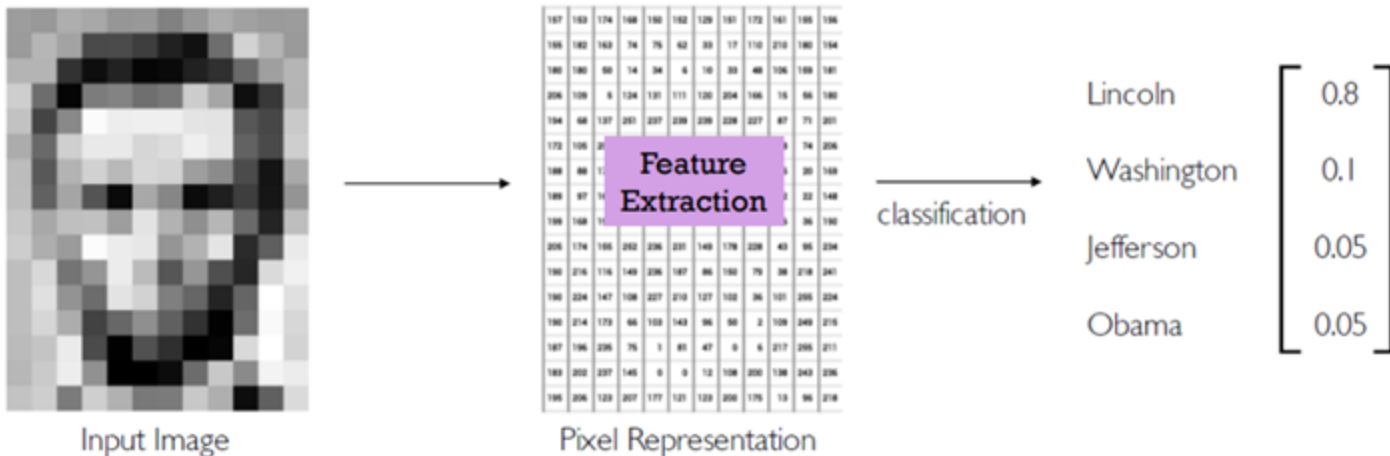
How robots' intelligent behaviors are guided by their interactive perception

How to close the loop!

# Quick review of DL

## Tasks in Computer Vision



- **Regression**: output variable takes continuous value
- **Classification**: output variable takes class label. Can produce probability of belonging to a particular class

# Quick review of DL

## High Level Feature Detection

Let's identify key features in each image category



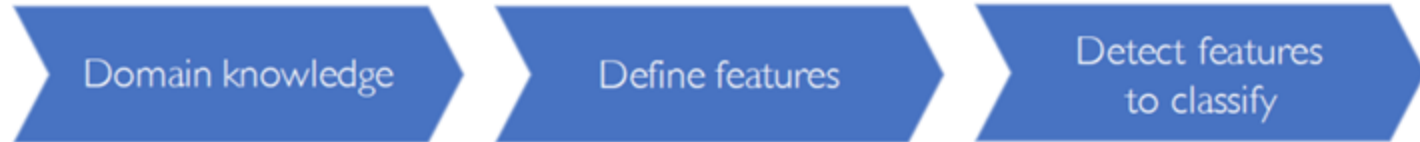Nose,
Eyes,
Mouth



Wheels,
License Plate,
Headlights



Door,
Windows,
Steps

# Quick review of DL

**Manual Feature Extraction**

Domain knowledge → Define features → Detect features to classify

Problems?

# Quick review of DL
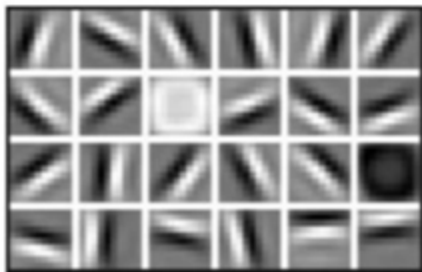
## Manual Feature Extraction

# Quick review of DL

Hand engineered features are time consuming, brittle and not scalable in practice

**Question 1** Can we learn the **underlying features** directly from data?

| Low Level Features | Mid Level Features | High Level Features |
| :---: | :---: | :---: |
|  |  |  |
| Lines & Edges | Eyes & Nose & Ears | Facial Structure |

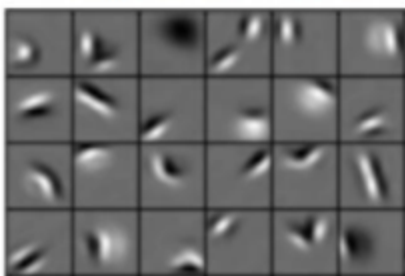# Quick review of DL

Can we **learn hierarchy of features** directly from the data instead of hand engineering?

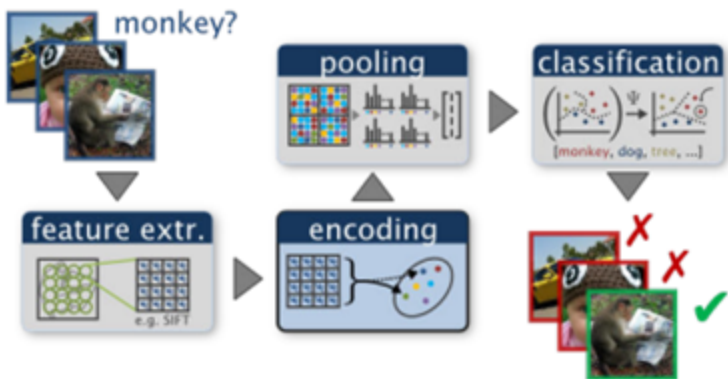| Low level features | Mid level features | High level features |
|---|---|---|
| Edges, dark spots | Eyes, ears, nose | Facial structure |



Slide Credit: Ava Soleimany, MIT

# Quick review of DL



Staged Visual Recognition Pipeline
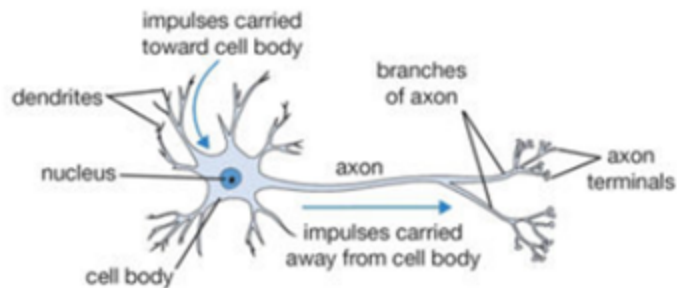
What is new since 1980s?
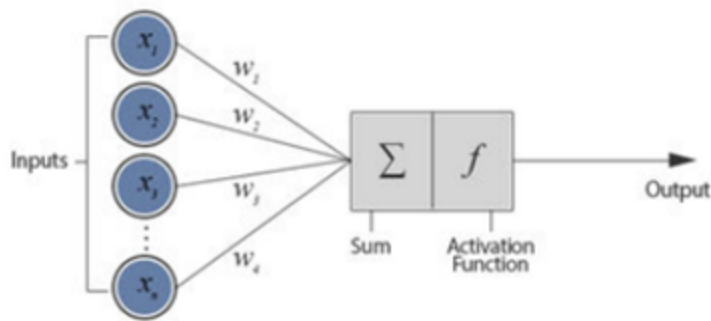
End-to-end Deep Learning

# Quick review of DL

## Biological Neuron versus Artificial Neural Network



**Biological Neuron**

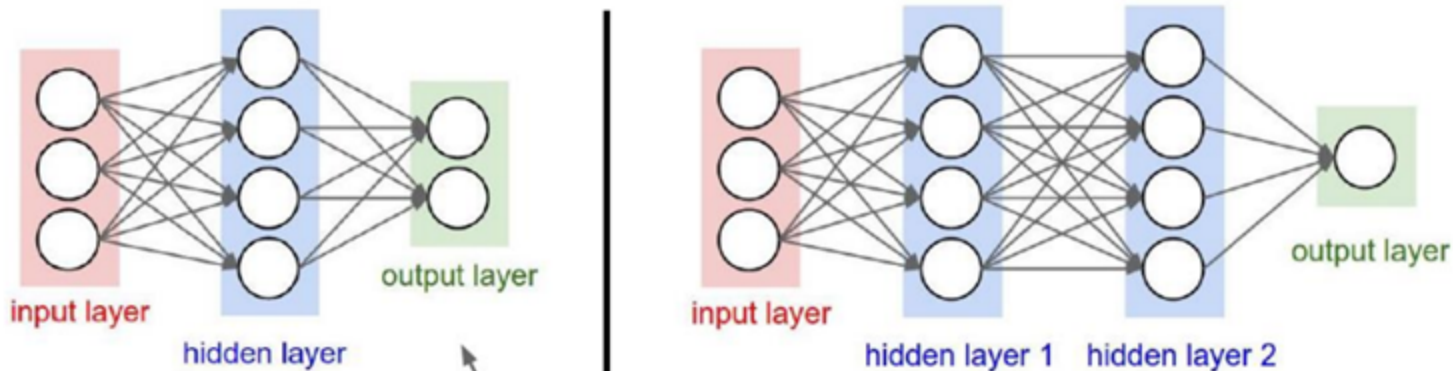Computational building block for the brain

**Artificial Neuron**

Computational building block for the neural network

**Note:** Many differences exist – be careful with the brain analogies!

[Dendritic Computation, Michael London and Michael Hausser 2015]

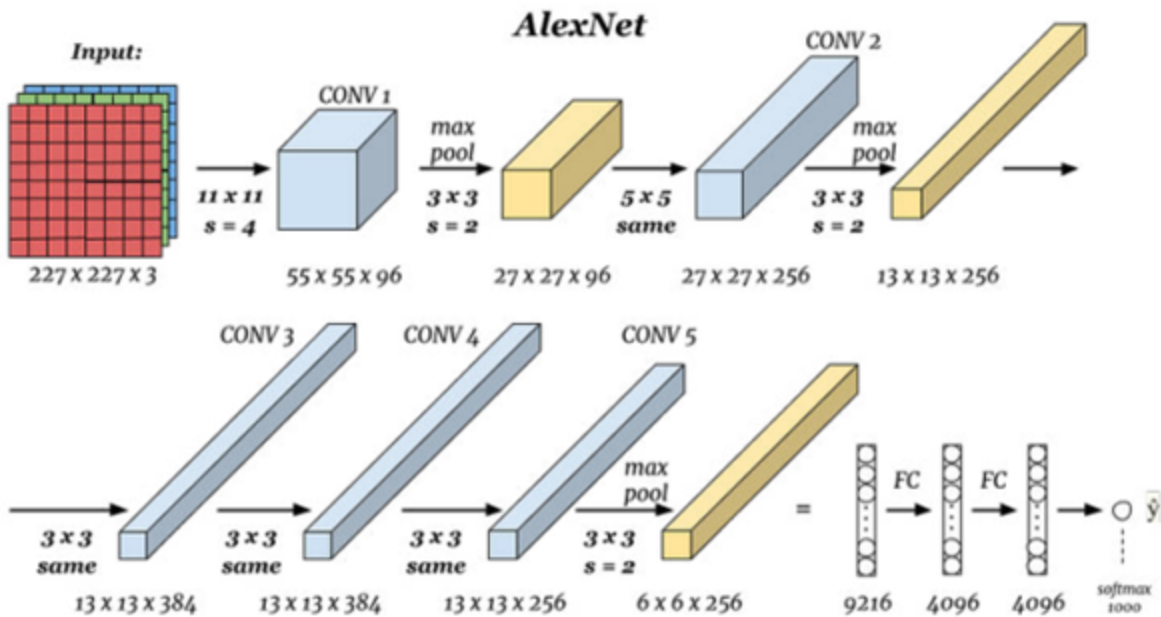# Quick review of DL

## Neural Networks: Architectures



"2-layer Neural Net", or
"1-hidden-layer Neural Net"

"Fully-connected" layers

"3-layer Neural Net", or
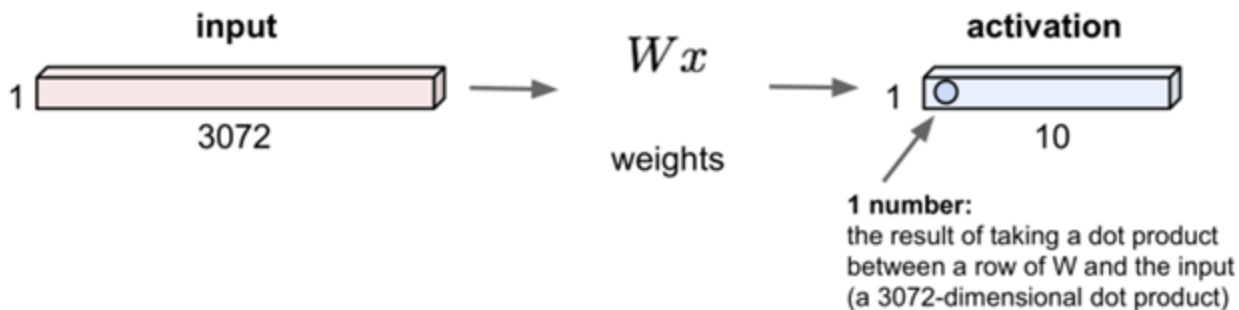"2-hidden-layer Neural Net"

# Quick review of DL



**AlexNet**

Input: 227 x 227 x 3 → CONV 1 (11 x 11, s = 4) → 55 x 55 x 96 → max pool (3 x 3, s = 2) → 27 x 27 x 96 → CONV 2 (5 x 5, same) → 27 x 27 x 256 → max pool (3 x 3, s = 2) → 13 x 13 x 256

CONV 3 (3 x 3, same) → 13 x 13 x 384 → CONV 4 (3 x 3, same) → 13 x 13 x 384 → CONV 5 (3 x 3, same) → 13 x 13 x 256 → max pool (3 x 3, s = 2) → 6 x 6 x 256 = 9216 → FC → 4096 → FC → 4096 → softmax 1000 → $\hat{y}$

https://indoml.com

# Quick review of DL

## Fully-Connected Layers

32x32x3 image -> stretch to 3072 x 1

**input**

1

3072

$Wx$

weights

**activation**

1

10

**1 number:**
the result of taking a dot product
between a row of W and the input
(a 3072-dimensional dot product)

What is the dimension of $W$?

[Source: Stanford CS231N]

# Quick review of DL

32x32x3 image -> preserve spatial structure
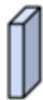


32 height

32 width

3 depth

# Quick review of DL


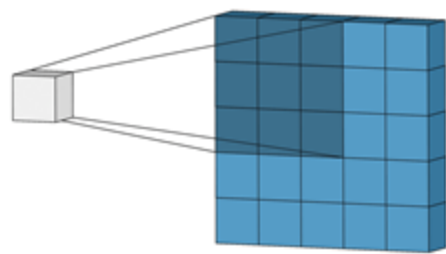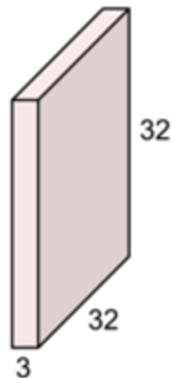
32x32x3 image

32

32

3

5x5x3 filter

**Convolve** the filter with the image
i.e. "slide over the image spatially,
computing dot products"

# Quick review of DL



32x32x3 image

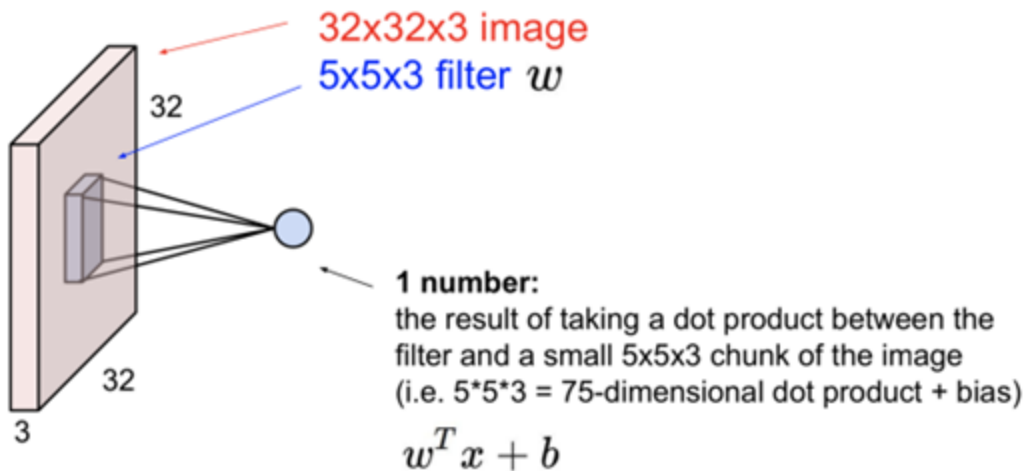Filters always extend the full depth of the input volume

32

32

3

5x5x3 filter

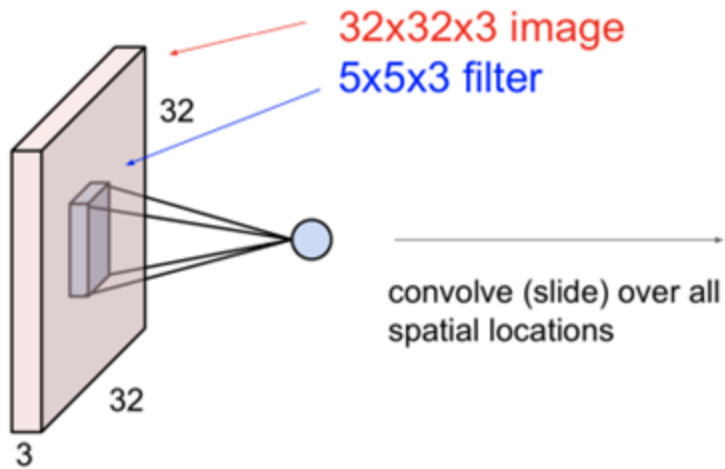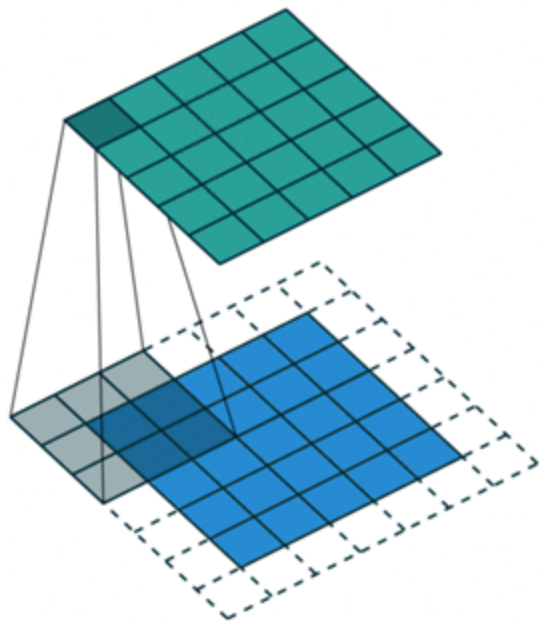**Convolve** the filter with the image i.e. "slide over the image spatially, computing dot products"
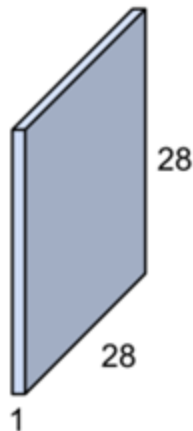
# Quick review of DL



32x32x3 image

5x5x3 filter $w$

32

32

3

**1 number:**
the result of taking a dot product between the
filter and a small 5x5x3 chunk of the image
(i.e. 5*5*3 = 75-dimensional dot product + bias)

$$w^T x + b$$

# Quick review of DL



32x32x3 image
5x5x3 filter

32

32

3

convolve (slide) over all spatial locations

**activation map**

28

28

1

# Quick review of DL

| | | | | |
|---|---|---|---|---|
| $3_0$ | $3_1$ | $2_2$ | 1 | 0 |
| $0_2$ | $0_2$ | $1_0$ | 3 | 1 |
| $3_0$ | $1_1$ | $2_2$ | 2 | 3 |
| 2 | 0 | 0 | 2 | 2 |
| 2 | 0 | 0 | 0 | 1 |

| | | |
|---|---|---|
| 12.0 | 12.0 | 17.0 |
| 10.0 | 17.0 | 19.0 |
| 9.0 | 6.0 | 14.0 |

consider a second, green filter

32x32x3 image
5x5x3 filter

32

32

3

convolve (slide) over all
spatial locations

activation maps

28

28

1

# Quick review of DL

For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:
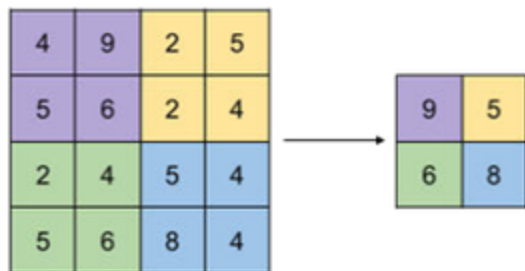


**activation maps**

Convolution Layer

32
32
3

28
28
6

We stack these up to get a "new image" of size 28x28x6!

# Quick review of DL



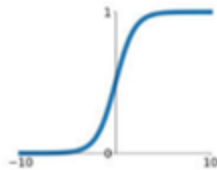Max Pooling

Avg Pooling

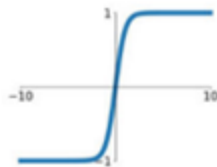https://indoml.com

# Quick review of DL



**Sigmoid**

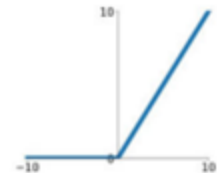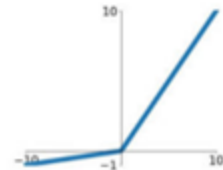$\sigma(x) = \frac{1}{1+e^{-x}}$

**tanh**
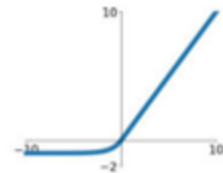
$\tanh(x)$

**ReLU**

$\max(0, x)$

**Leaky ReLU**

$\max(0.1x, x)$

**Maxout**

$\max(w_1^T x + b_1, w_2^T x + b_2)$

**ELU**

$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$
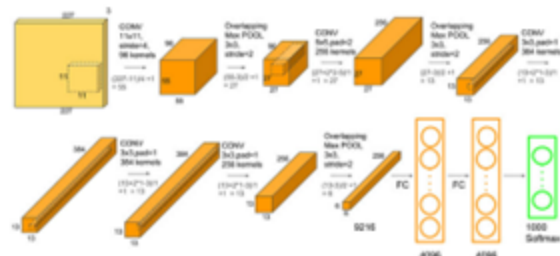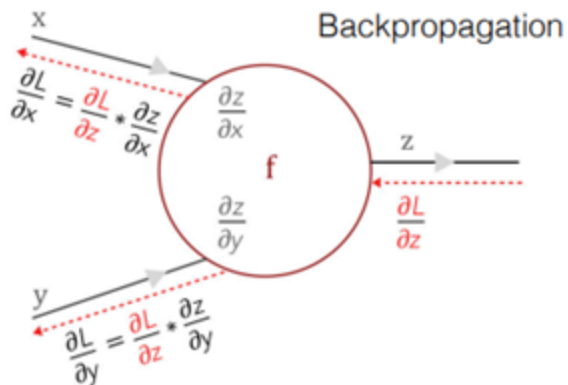
# Quick review of DL



LeNet

ResNet

VGG-16

AlexNet

# Quick review of DL

# Quick review of DL

Tutorial coming in late September / early October

```python
import torch
from torch import nn

class MNISTClassifier(nn.Module):

    def __init__(self):
        super(MNISTClassifier, self).__init__()

        # mnist images are (1, 28, 28) (channels, width, heig
        self.layer_1 = torch.nn.Linear(28 * 28, 128)
        self.layer_2 = torch.nn.Linear(128, 256)
        self.layer_3 = torch.nn.Linear(256, 10)

    def forward(self, x):
        batch_size, channels, width, height = x.size()

        # (b, 1, 28, 28) -> (b, 1*28*28)
        x = x.view(batch_size, -1)

        # layer 1
        x = self.layer_1(x)
        x = torch.relu(x)

        # layer 2
        x = self.layer_2(x)
        x = torch.relu(x)

        # layer 3
        x = self.layer_3(x)

        # probability distribution over labels
        x = torch.log_softmax(x, dim=1)

        return x
```

# Quick review of DL

### Online Courses

- CS231N: Convolutional Neural Networks for Visual Recognition

  http://cs231n.stanford.edu/

- MIT 6.S191: Introduction to Deep Learning

  http://introtodeeplearning.com/

### Textbooks:

- Deep Learning. Ian Goodfellow, Yoshua Bengio, Aaron Courville

  http://www.deeplearningbook.org/



DEEP LEARNING
Ian Goodfellow, Yoshua Bengio, and Aaron Courville