

# **Visionary Course – Energy AI**

## **Week 10**

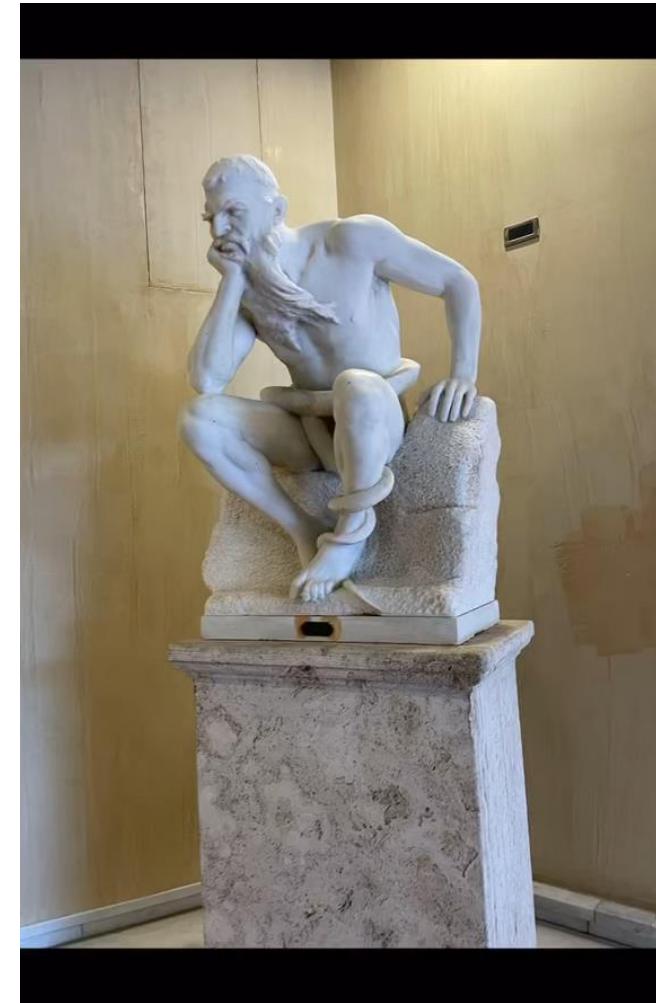
**Seokju Lee**

# Semantic Segmentation



# Photo Segmentation on My iPhone

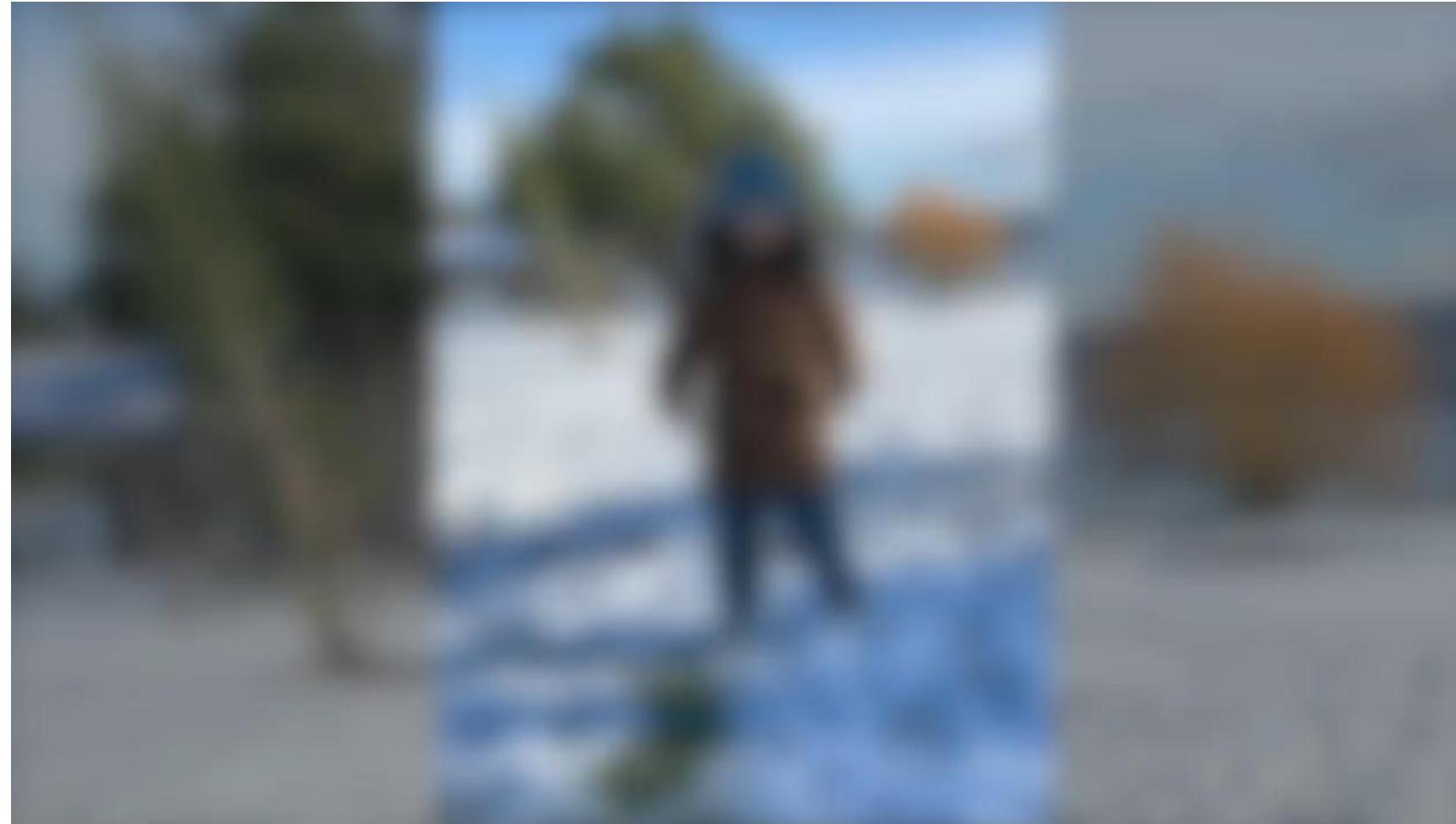
In the newest version of **iPhone** (iOS 16),



**Image segmentation!**

# Image Segmentation + 3D Computer Vision

**3D Photography** is available,



# Computer Vision Tasks

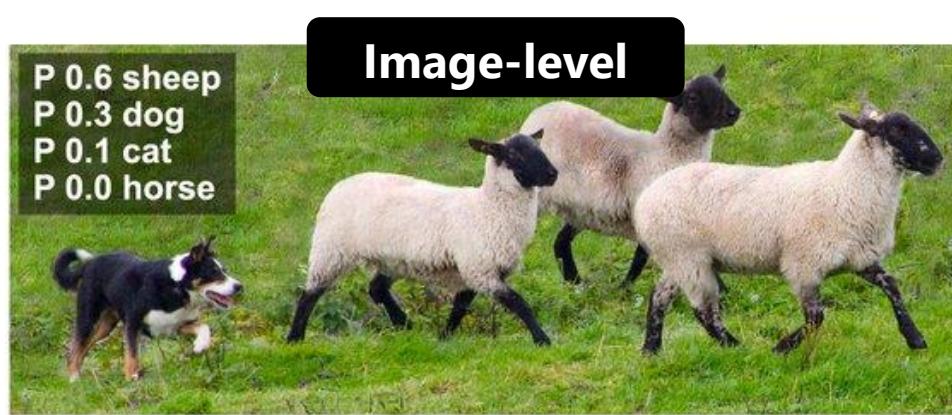
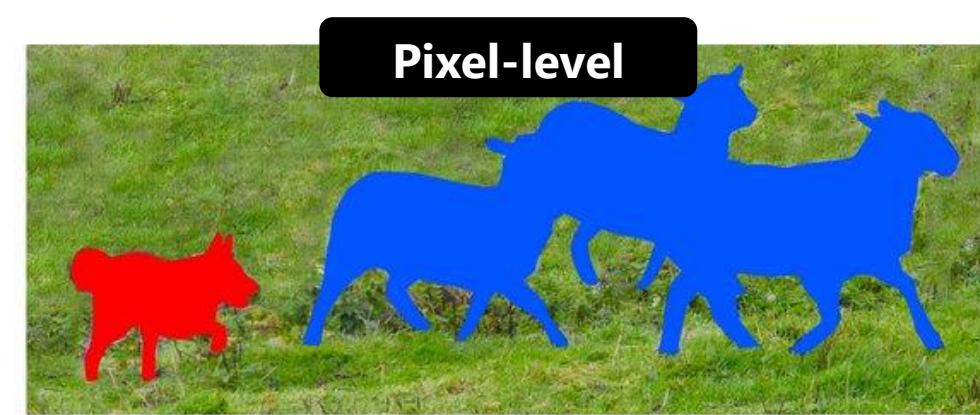
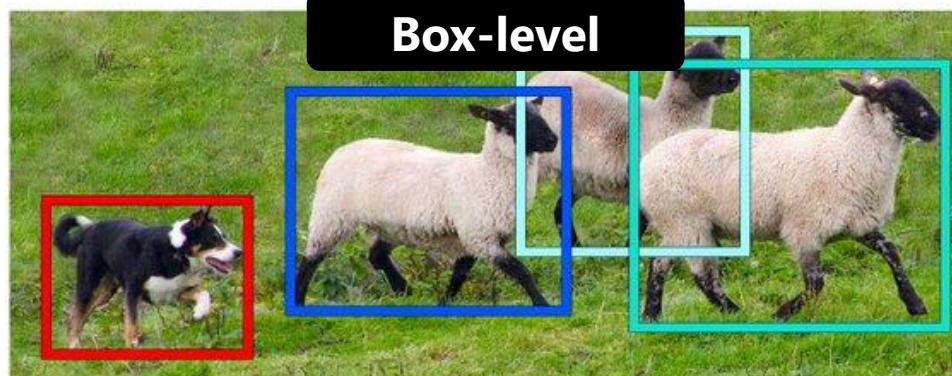


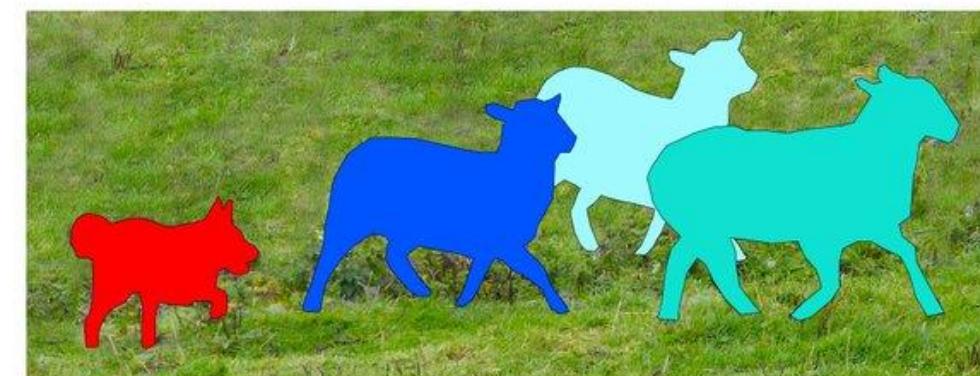
Image Recognition



Semantic Segmentation



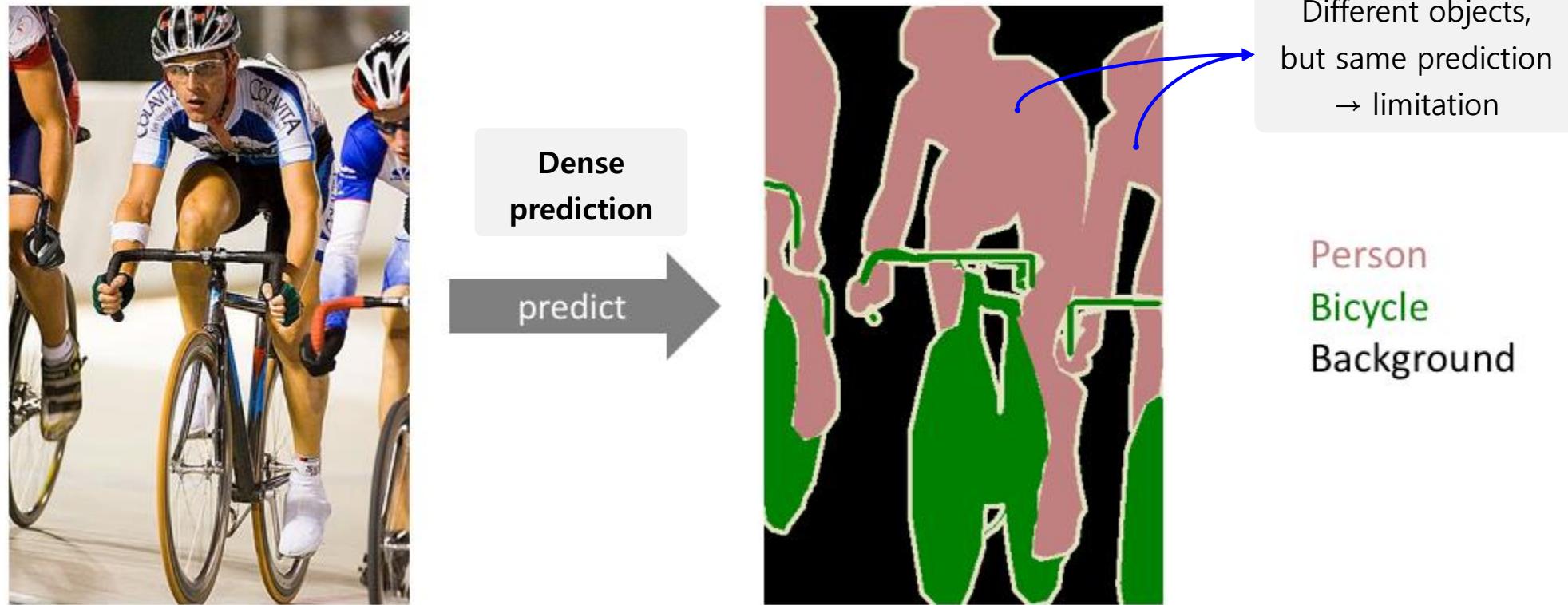
Object Detection



Instance Segmentation

# Semantic Segmentation

: Task of assigning **a label** for **every pixel** of an image with a corresponding **class**



# Semantic Segmentation: Input & Output



segmented

- 1: Person
- 2: Purse
- 3: Plants/Grass
- 4: Sidewalk
- 5: Building/Structures

Ground Truth (GT) → supervised learning

3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	5	5	5	5	5	
3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	5	5	5	5	5
3	3	3	3	3	3	3	3	3	3	1	1	1	3	3	3	3	3	3	3	5	5	5	5	5	5
3	3	3	3	3	3	3	3	3	3	1	1	1	1	1	3	3	3	3	3	5	5	5	5	5	5
3	3	3	3	3	3	3	3	3	3	1	1	1	1	3	3	3	3	3	3	5	5	5	5	5	5
5	5	3	3	3	3	3	3	3	3	1	1	1	3	3	3	3	3	3	5	5	5	5	5	5	5
4	4	3	4	1	1	1	1	1	1	1	1	1	1	4	4	4	4	4	4	5	5	5	5	5	5
4	4	3	4	1	1	1	1	1	1	1	1	1	1	4	4	4	4	4	4	5	5	5	5	5	5
4	4	4	1	1	1	1	1	1	1	1	1	1	1	1	4	4	4	4	4	4	4	4	4	4	4
3	3	3	1	1	1	1	1	1	1	1	1	1	1	1	4	4	4	4	4	4	4	4	4	4	4
3	3	3	1	2	2	1	1	1	1	1	1	1	1	1	4	4	4	4	4	4	4	4	4	4	4

Input

→ Input image resolution:  $W \times H \times 3$

Semantic Labels

Note that this is a low-resolution prediction map for visual clarity.

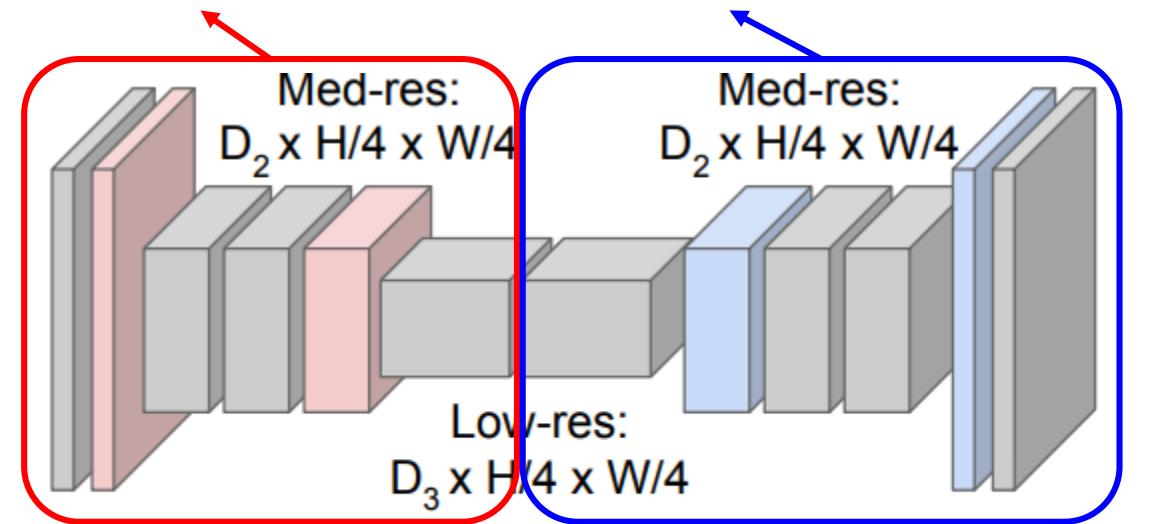
In reality, the segmentation label resolution should match the original input's resolution ( $W \times H \times 3$ ).

# Fully Convolutional Approach

Design a network as a bunch of convolutional layers,  
with **downsampling** and **upsampling** inside the network!



Input:  
 $3 \times H \times W$



Predictions:  
 $H \times W$

**Downsampling:** Feature extraction by pooling, strided convolution (encoder)

**Upsampling:** Recover of spatial resolution by unpooling, strided transpose convolution (decoder)

# Semantic Segmentation on Jetson Nano

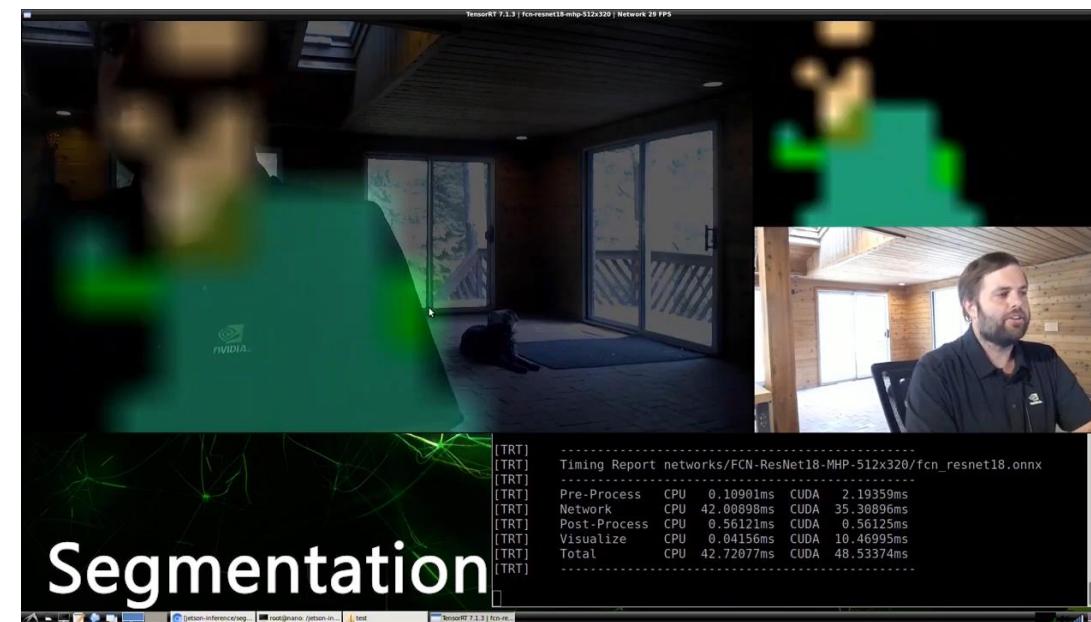
# Experiments – Semantic Segmentation

### Live semantic segmentation with visualization ###

\*Your basic workspace is here: “[cd ~/jetson-inference/build/aarch64/bin](#)”

Q1. Run “[python segnet.py --flip-method=rotate-180](#)”. What is on the screen? What is different from detectnet in the previous class (object detection)? Could you guess the meaning of the different colors on the screen?

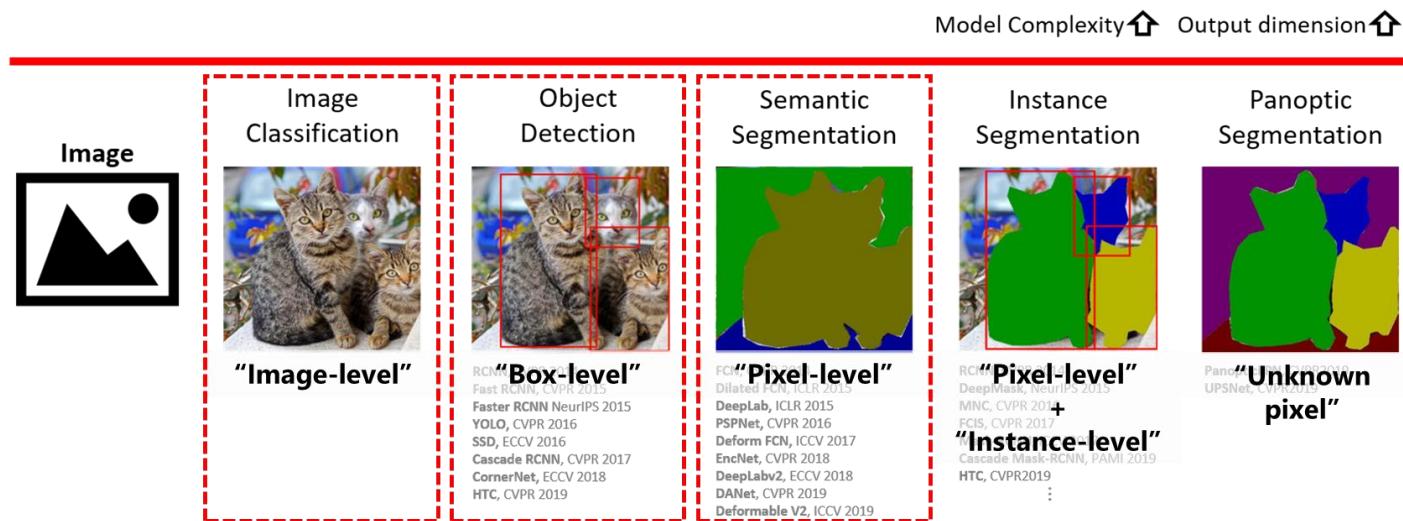
Q2. Go to the linked page (<https://github.com/dusty-nv/jetson-inference/blob/master/docs/segnet-console-2.md>). Please read the list of segmentation models. Try to use different models, trained on different datasets (upto your preference). For example, run “[python segnet.py --network=fcn-resnet18-cityscapes-1024x512 --flip-method=rotate-180](#)”, and let the camera see the driving video.



# Summary

## Basic computer vision tasks (semantic tasks)

- Image classification (AlexNet, VGG, GoogleNet, ResNet): Basic deep neural networks
- Object detection (R-CNN, SSD): Trade-off between two-stage & one-stage detectors
- Semantic segmentation (Fully Convolutional Network): Encoder (feature extraction) + Decoder (upsampling)



→ Combinatorial works of object  
detection and semantic segmentation

→ We will move on to **geometric tasks!**

# Q&A

**KENTECH**  
Korea Institute of Energy Technology

# 3D Visual Perception

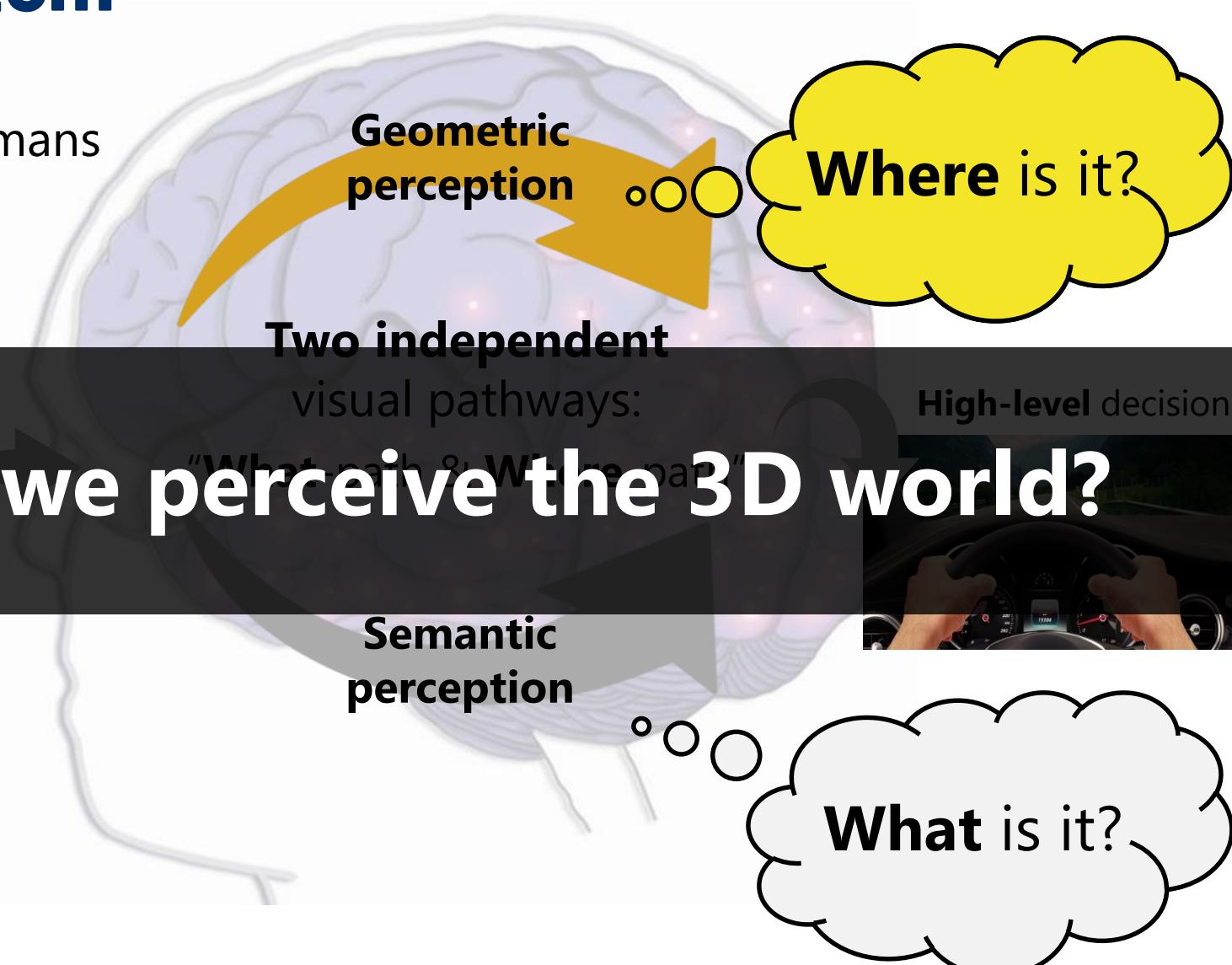
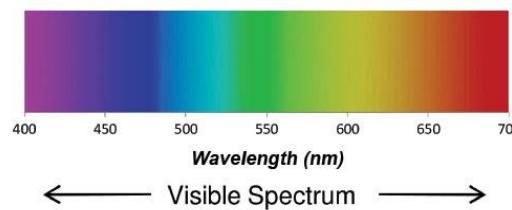


# Human Visual System

"About **half** of neocortex in humans is devoted to **vision**." [1]



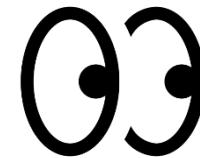
**Low-level visual signal**



[1] Barton, Robert A. "Visual specialization and brain evolution in primates." *Proceedings of the Royal Society of London* (1998).

[2] M. A. Goodale, et al., "Separate visual pathways for perception and action." *Trends in Neurosciences* (1992).

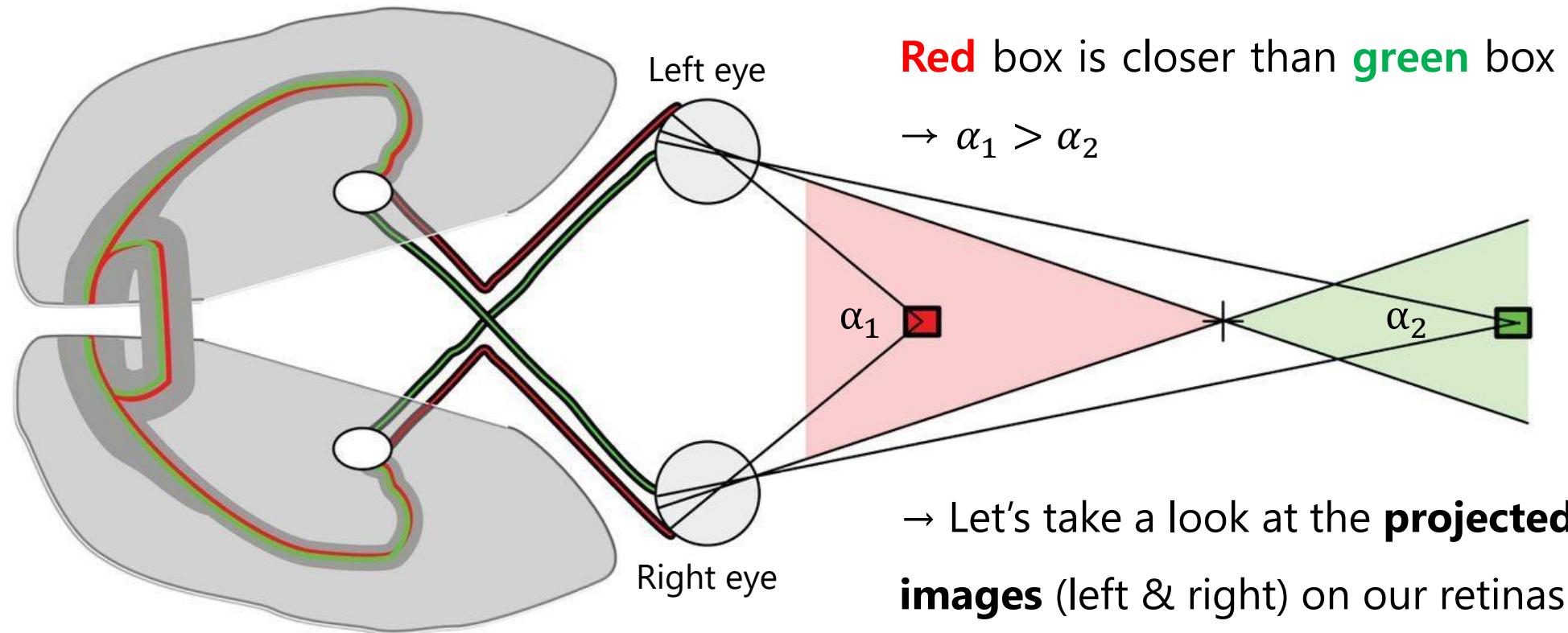
# Human 3D Perception: Binocular System



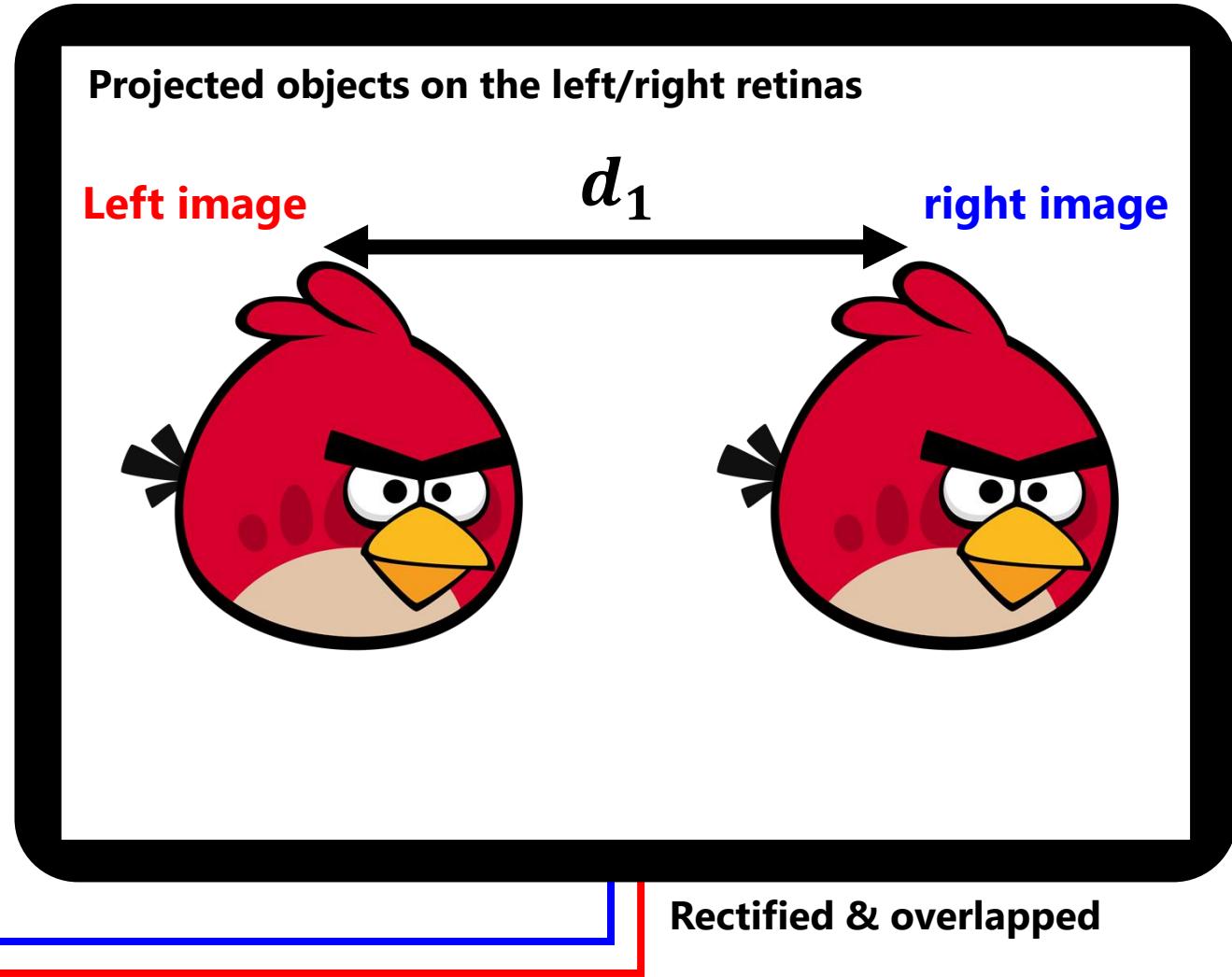
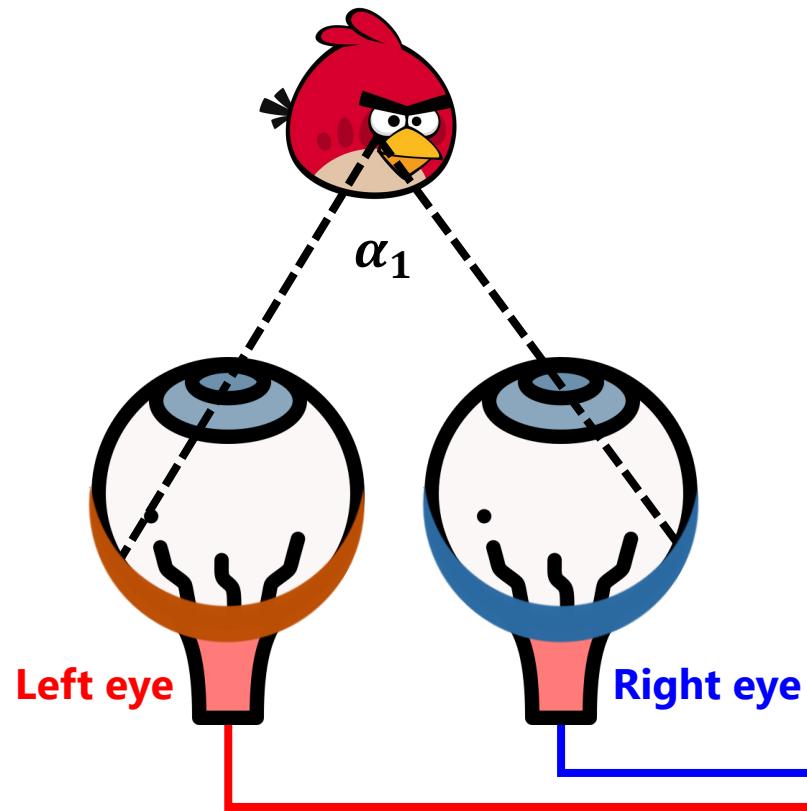
## Binocular correspondence

Two eyes

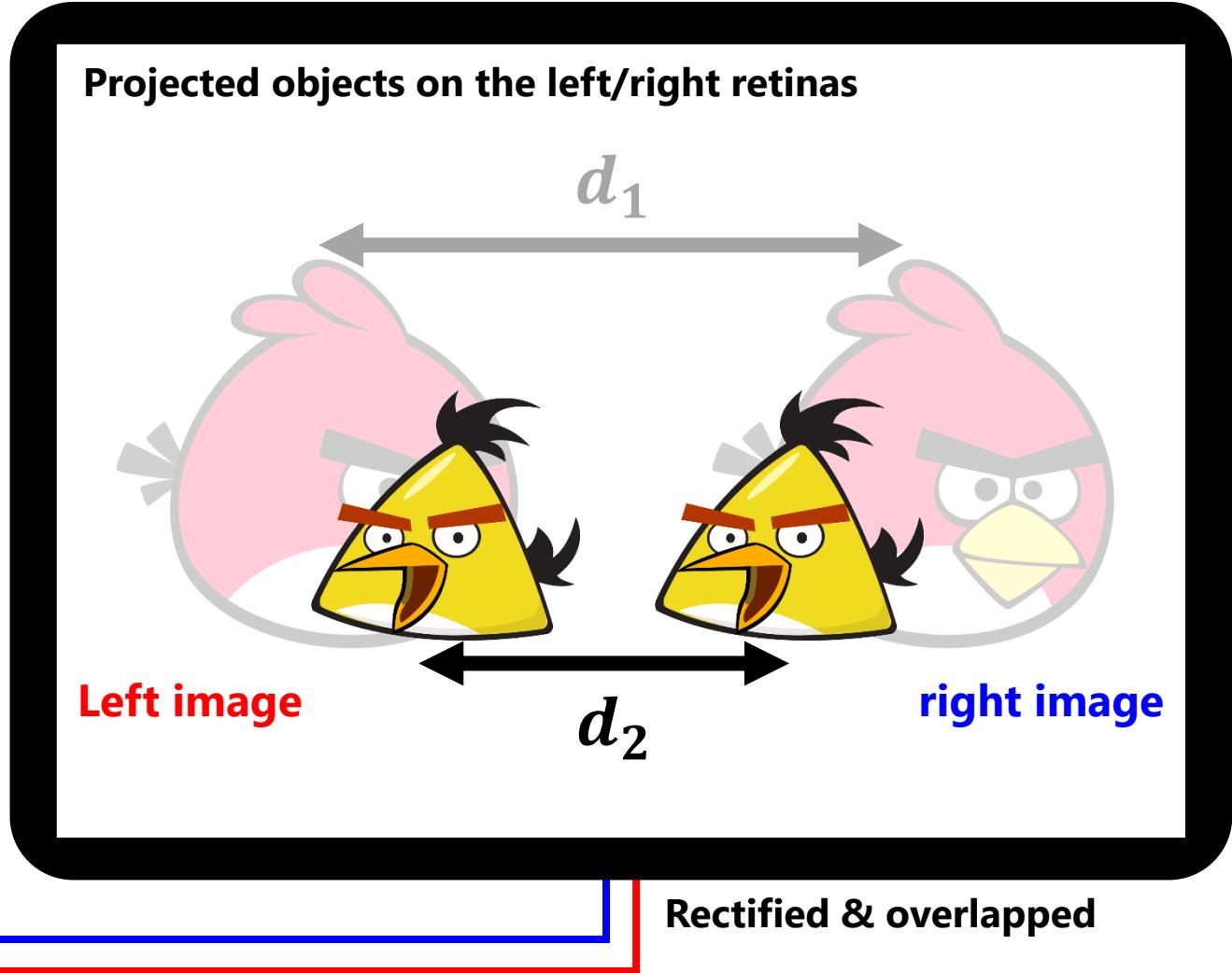
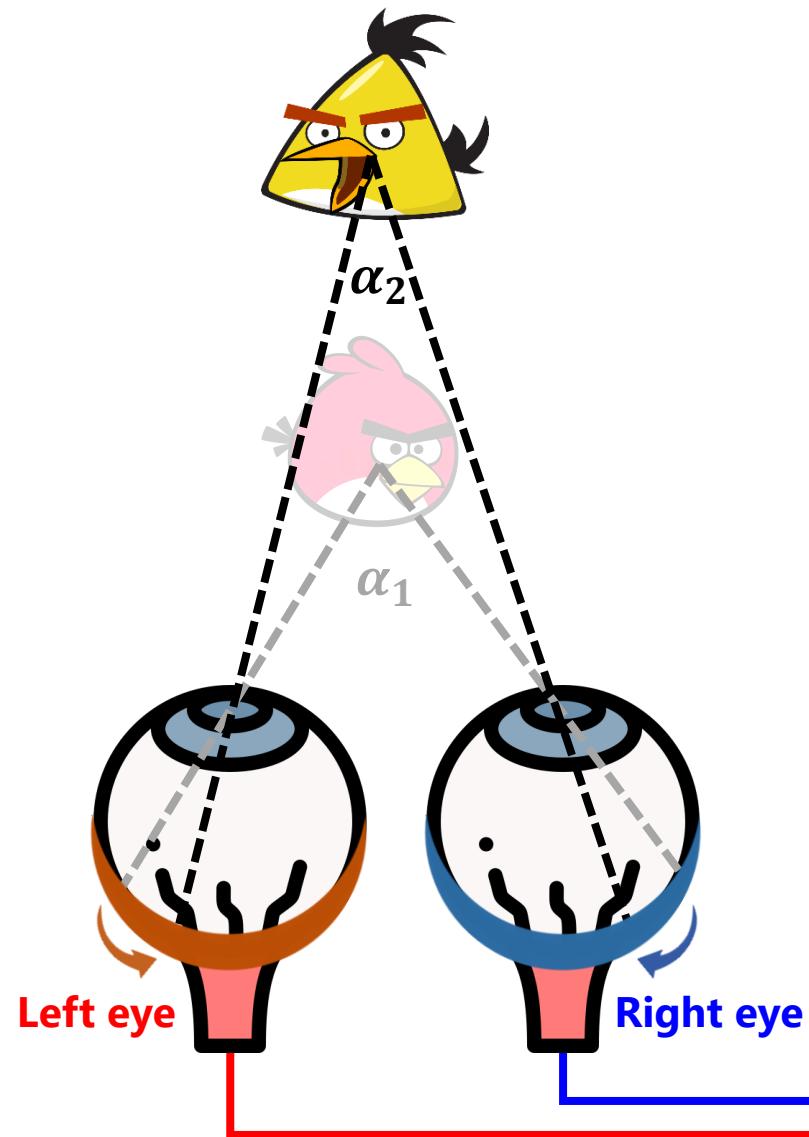
Finding matching points



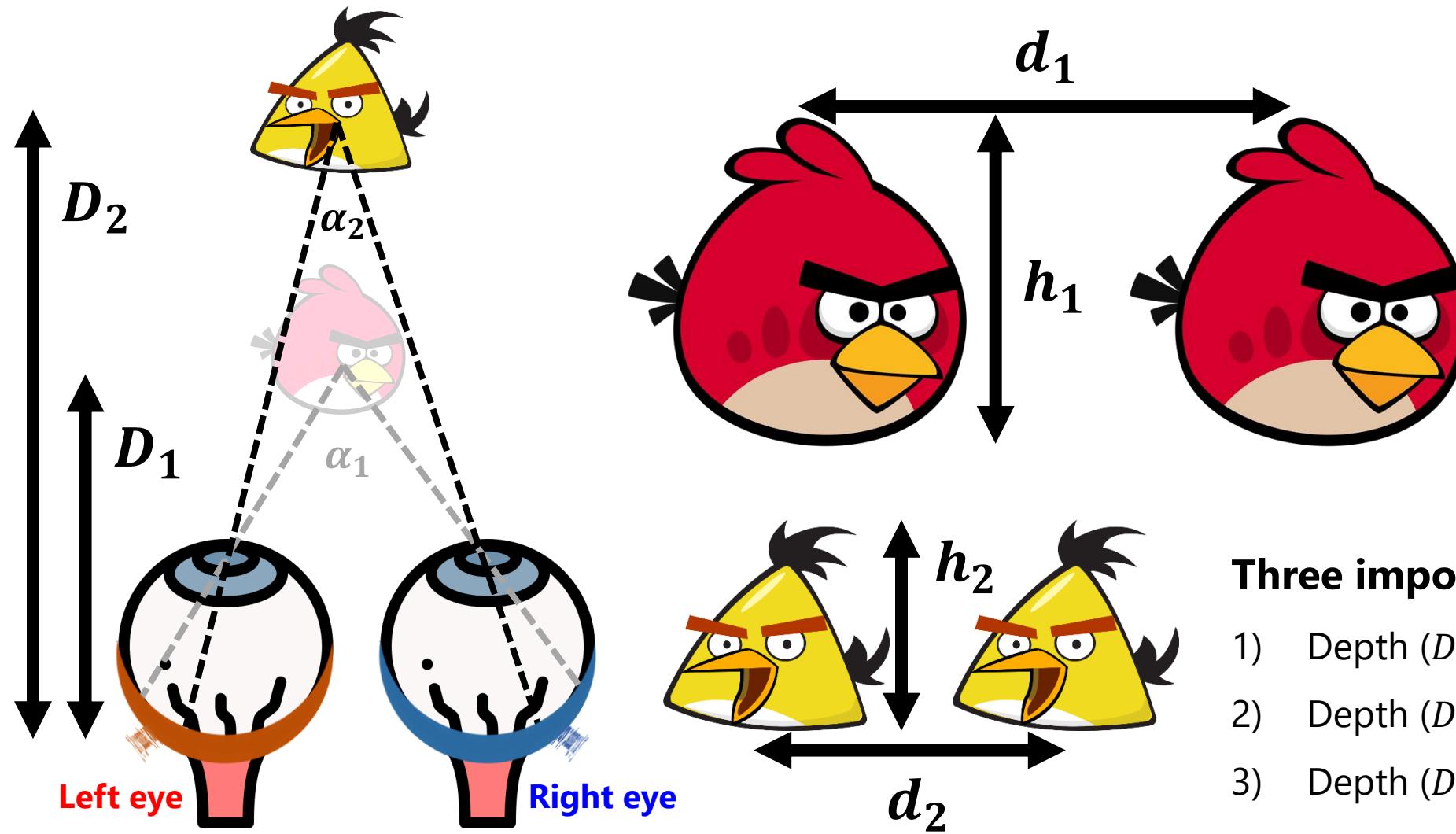
# Binocular Correspondence: Near Objects



# Binocular Correspondence: Far Objects



# Binocular Correspondence: Disparity ( $d$ )



## Three important facts

- 1) Depth ( $D$ )  $\uparrow \rightarrow$  Angle ( $\alpha$ )  $\downarrow$
- 2) Depth ( $D$ )  $\uparrow \rightarrow$  Disparity ( $d$ )  $\downarrow$
- 3) Depth ( $D$ )  $\uparrow \rightarrow$  Object's size ( $h$ )  $\downarrow$

# Human perception is remarkably flexible!

## Other depth cues to perceive 3D world?

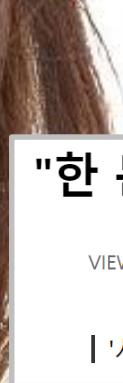
# Perceiving 3D World with Monocular System



## Driving With One Eye

*Department of Transportation, USA*

Not only is it possible to drive with one eye (assuming that you have good vision in your remaining eye) it's also **legal** in many states. Though there isn't a federal law dictating whether people with monocular vision can drive, it is up to each state to determine these regulations.<sup>[7]</sup> Of course, like everything else, it will take some time and practice to get used to driving with one eye, so you may want to look into specialized driving classes in your area.



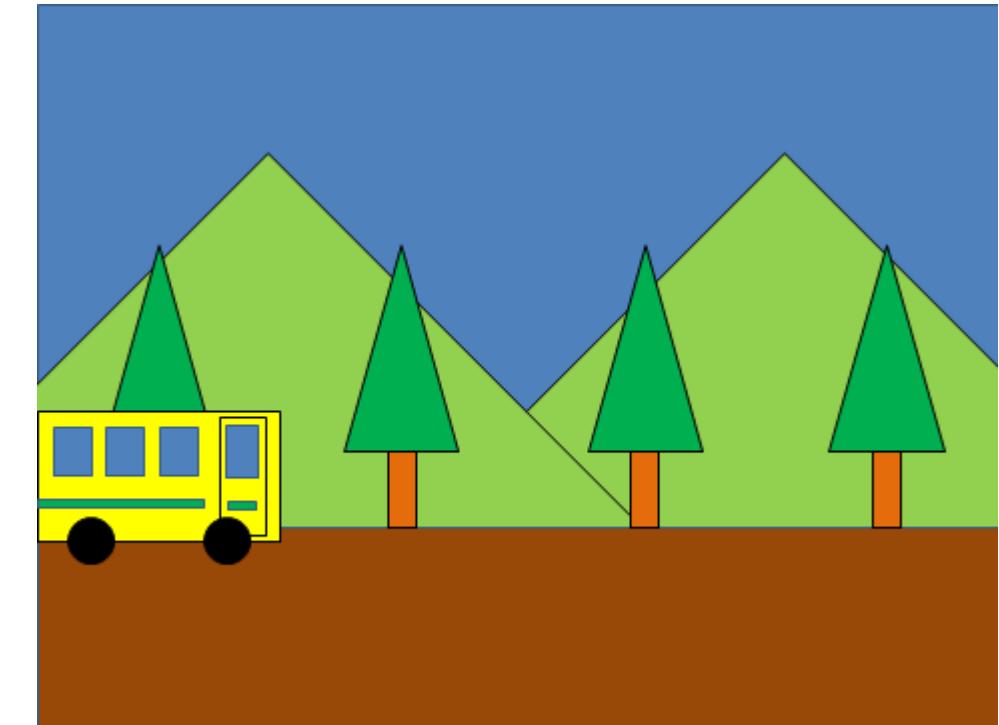
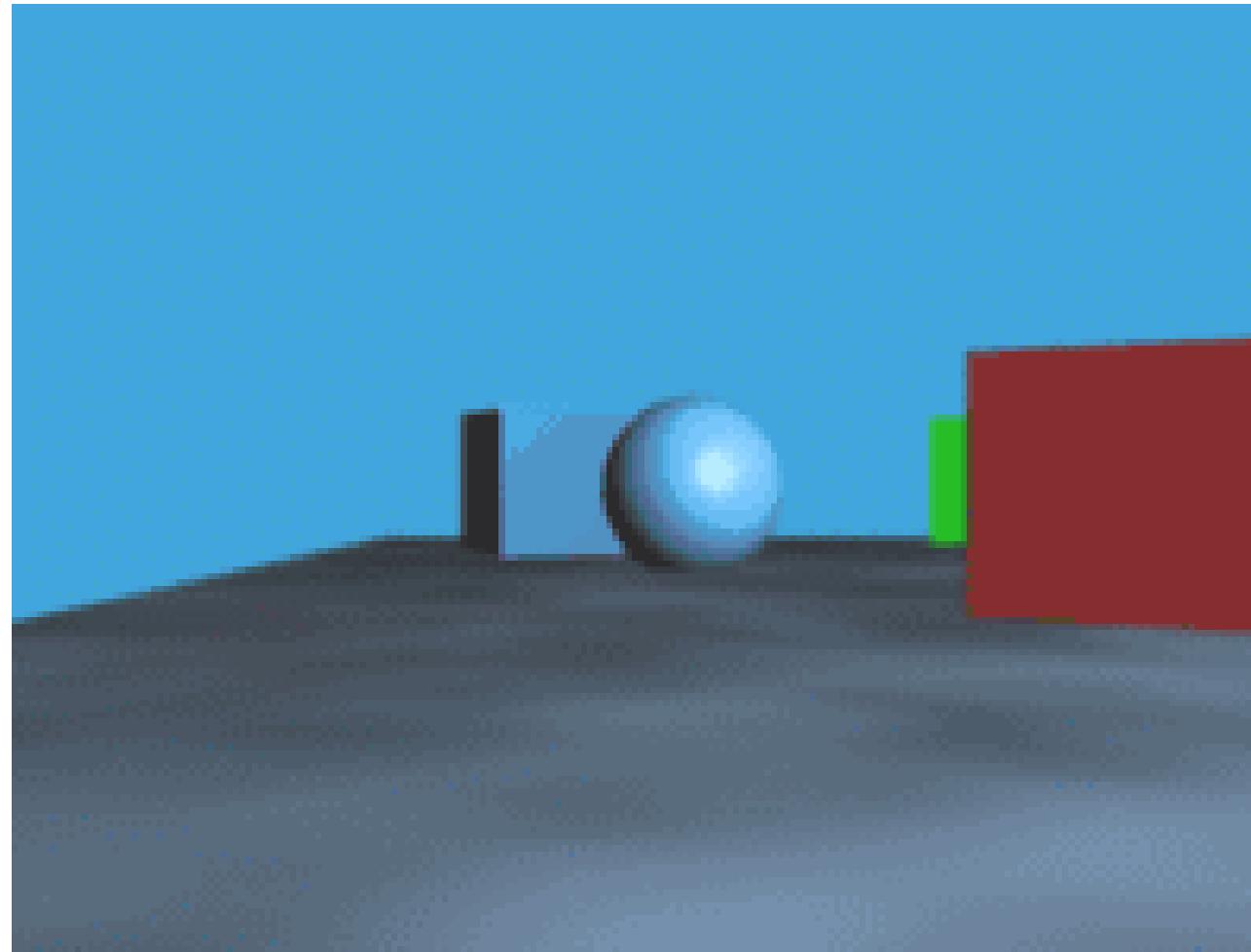
"한 눈 잘 보이면" 시각장애인도 '1종 보통' 면허 땐다

VIEW 19,534 | 2016.04.18 05:20



| '시력 0.8·일정 수준 이상 시야' 경찰, 단안 시각장애 면허기준 마련

# Monocular Depth Cues (1): Motion



## “Motion parallax”

→ Differences in image **motions** between objects at different **depths** (motion  $\uparrow \rightarrow$  depth  $\downarrow$ )

# Motion Parallax in Film Industry

## Walt Disney's Multiplane Camera



→ Multiplane images allow us to feel the 3D effect !

# Motion Parallax by Animals

Strange behaviors of animals...



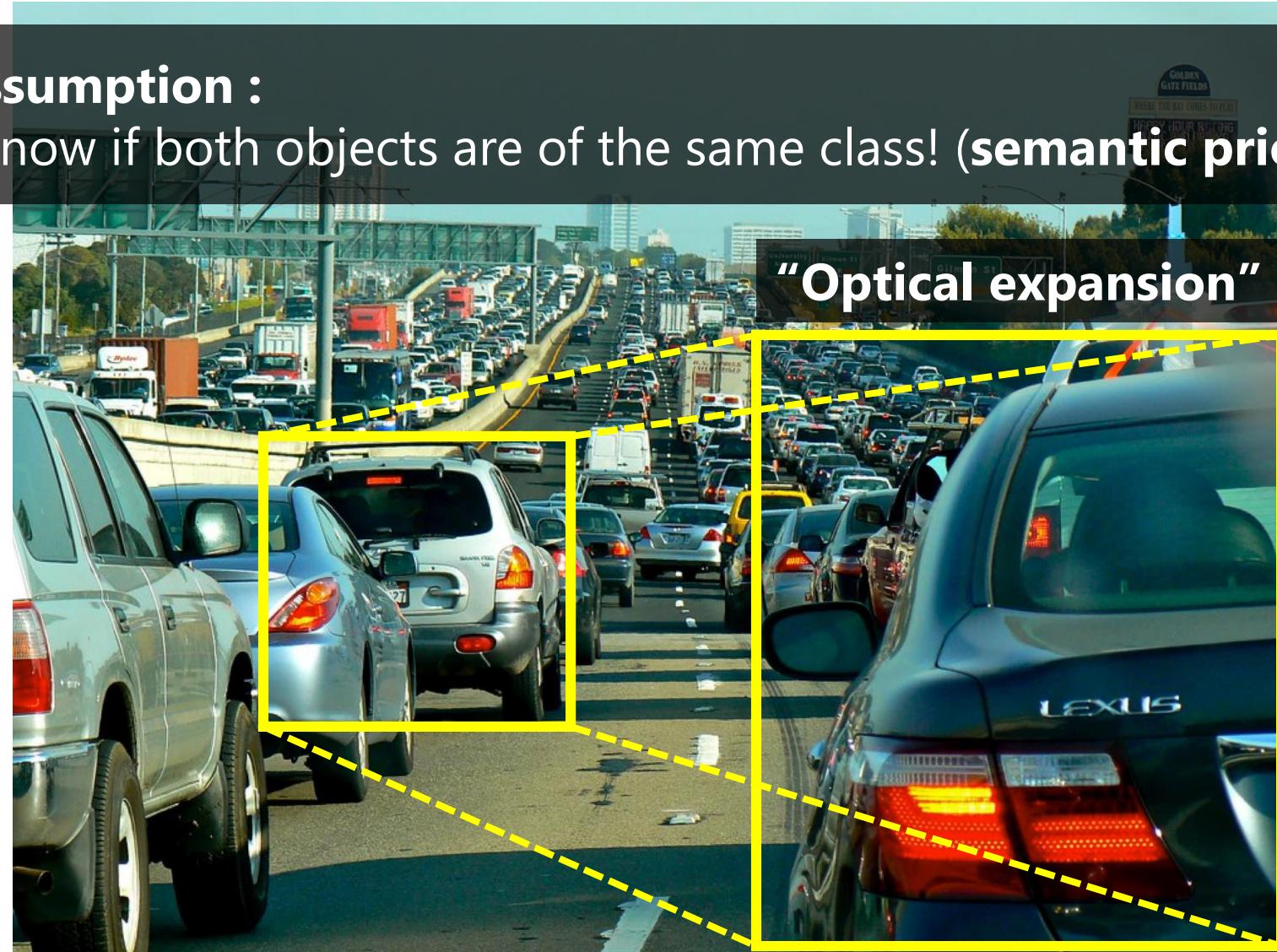
- [1] YouTube, "Getting down with some biggie!"
- [2] YouTube, "Super cute budgie bobbing his head"
- [2] YouTube, "Headbanging/bobbing Turtle"

Have a reason... 3D perception!

# Monocular Depth Cues (2): Relative Size

Here is an assumption :

We need to know if both objects are of the same class! (**semantic prior** knowledge)

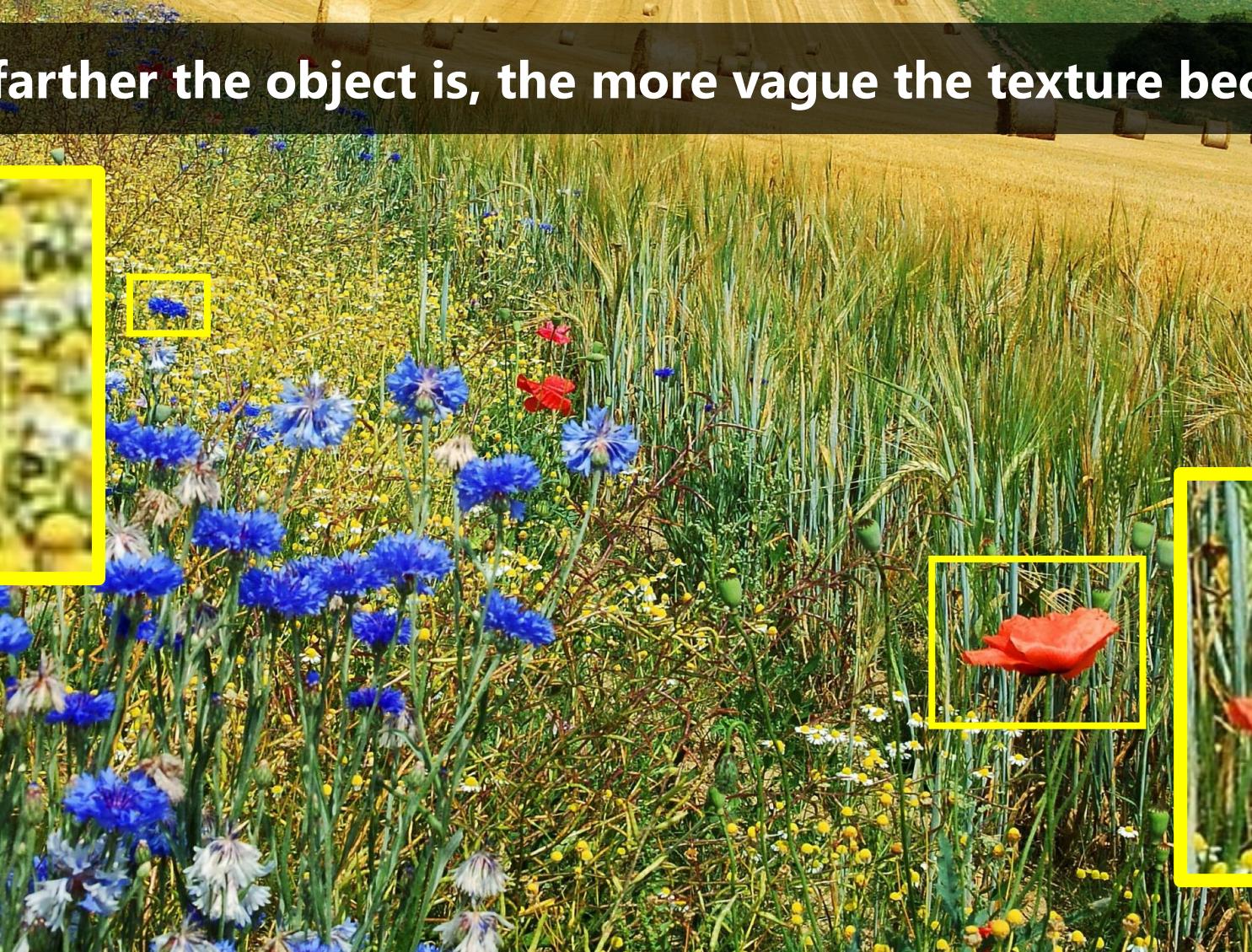


# Monocular Depth Cues (3): Texture Gradient

**"The farther the object is, the more vague the texture becomes"**

$$\left| \frac{\partial I}{\partial x} \right|, \left| \frac{\partial I}{\partial y} \right| \downarrow$$

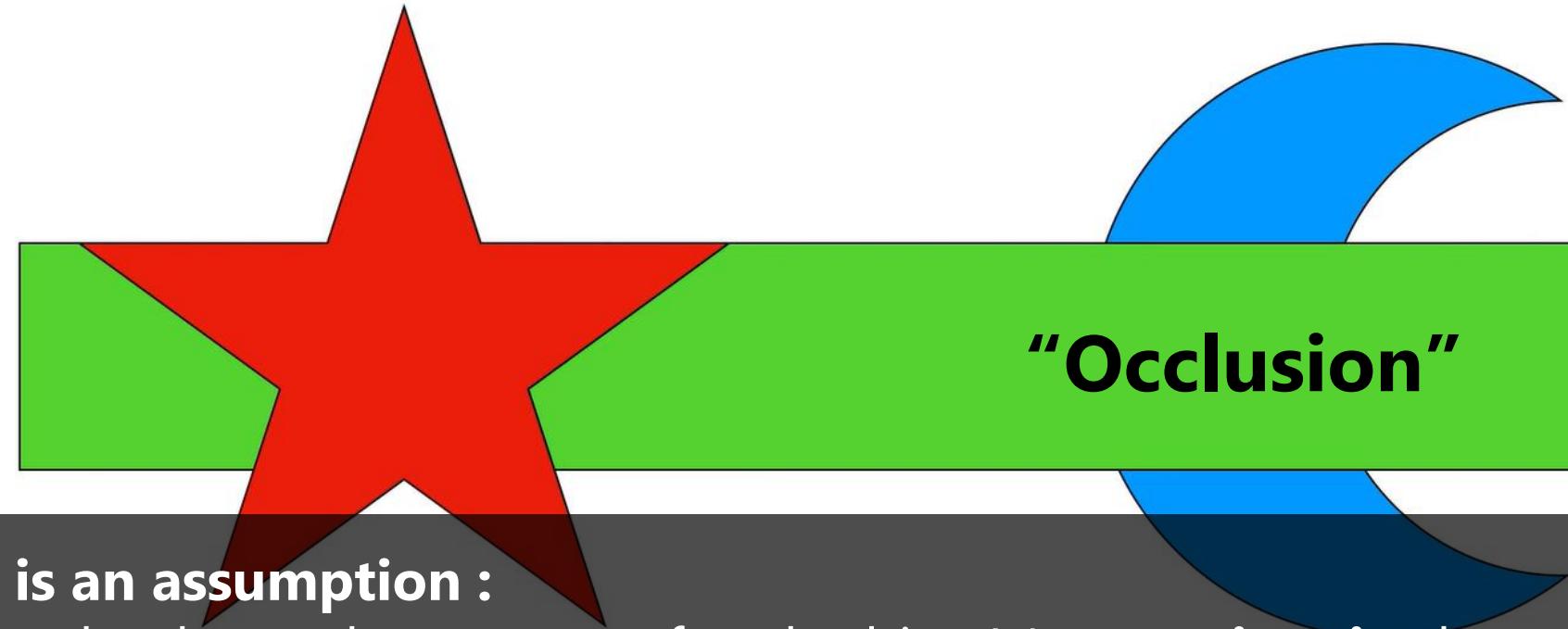
$$\left| \frac{\partial I}{\partial x} \right|, \left| \frac{\partial I}{\partial y} \right| \uparrow$$



# Monocular Depth Cues (4): Interposition

**Object A** is blocking our view of **Object B**

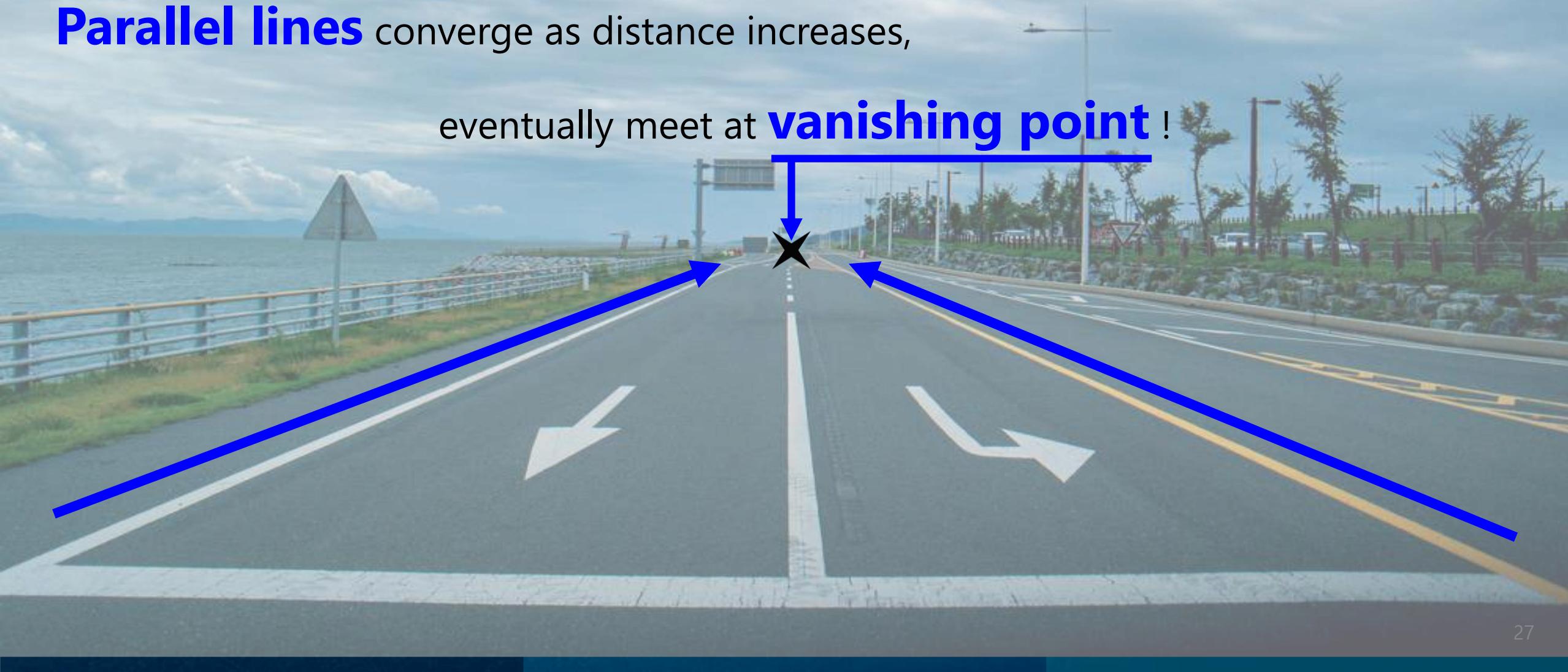
→ **Object A** is closer to us than **Object B**



# Monocular Depth Cues (5): Linear Perspective

**Parallel lines** converge as distance increases,

eventually meet at **vanishing point** !



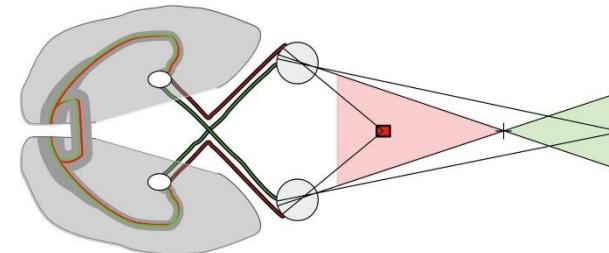
# Optical Illusion by False Linear Perspective



# Summary of 3D Perception: from Human to Machine

## Binocular correspondence

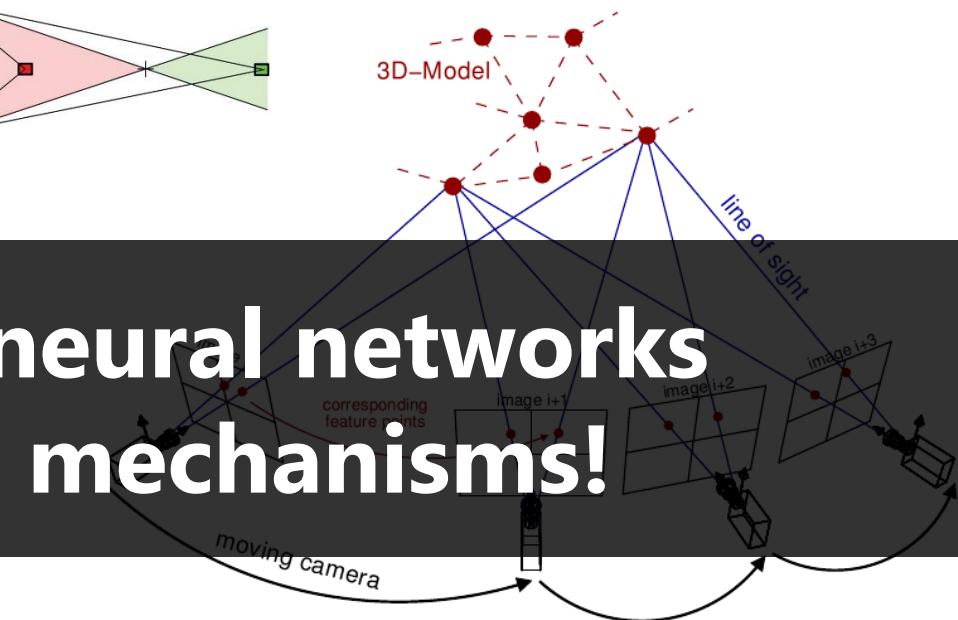
→ "Stereo matching" in computer vision



## Motion parallax

→ "Structure-from-Motion (SfM)" in computer vision

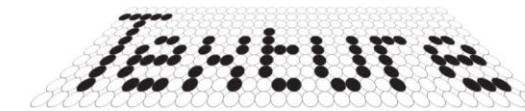
→ **Applied to train deep neural networks with different learning mechanisms!**



## Other monocular cues

shadow

SIZE  
occlusion



BLUR

Shading

perspective

# 3D Computer Vision Applications

3D Reconstruction



Speed 4x

CoEx, IROS'21



**BLOCK-NERF**

RESULTS

Matthew Tancik<sup>1\*</sup>   Vincent Casser<sup>2</sup>   Xincheng Yan<sup>2</sup>  
Sabeek Pradhan<sup>2</sup>   Ben Mildenhall<sup>3</sup>   Pratul Srinivasan<sup>3</sup>  
Jonathan T. Barron<sup>3</sup>   Henrik Kretzschmar<sup>2</sup>

<sup>1</sup>UC Berkeley   <sup>2</sup>Waymo   <sup>3</sup>Google Research

\*Work done as an intern at Waymo

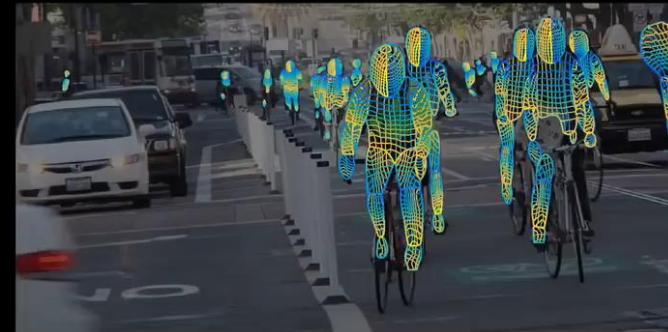
**Google**

Waymo's AI Recreates San Francisco From 2.8 Million Photos! 🚗, 2022

A screenshot of a 3D reconstruction application showing a street scene with a car and a building. The image includes text overlays for authors and institutions.

## DensePose:

Dense Human Pose Estimation In The Wild



Riza Alp Güler \*   Natalia Neverova   Iasonas Kokkinos  
*INRIA, CentraleSupélec*   *Facebook AI Research*   *Facebook AI Research*

\* Riza Alp Güler was with Facebook AI Research during this work.

DensePose, CVPR'18

3D human pose estimation



EG3D, 2021

3D face reconstruction

# Q&A

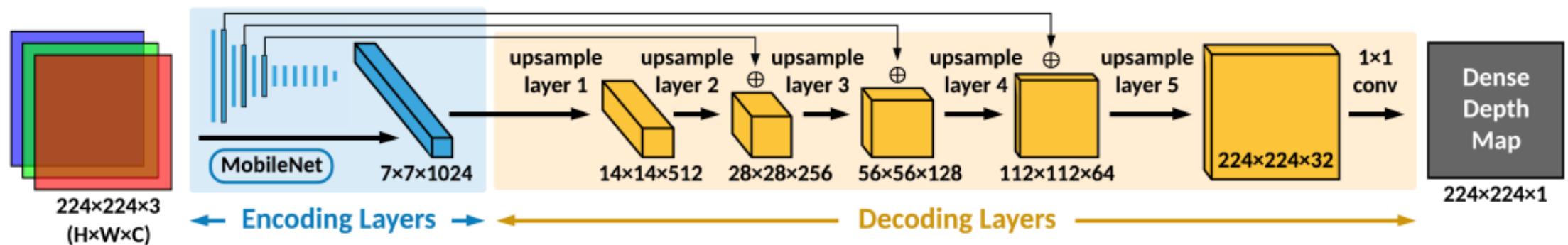
**KENTECH**  
Korea Institute of Energy Technology

(1) (2)

# Depth & Pose Estimation on Jetson Nano

# Reference – Depth Estimation

Monocular depth estimation (<https://github.com/dusty-nv/jetson-inference/blob/master/docs/depthnet.md>)

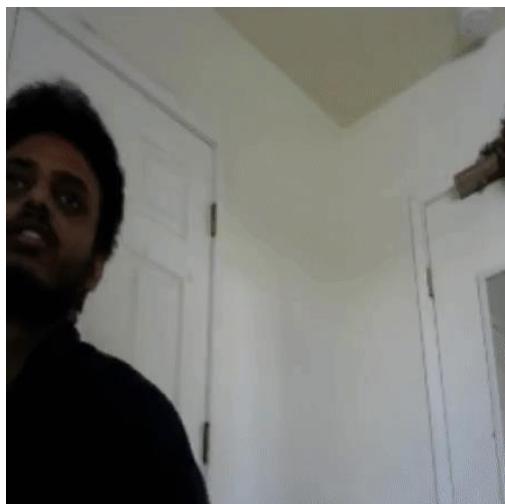
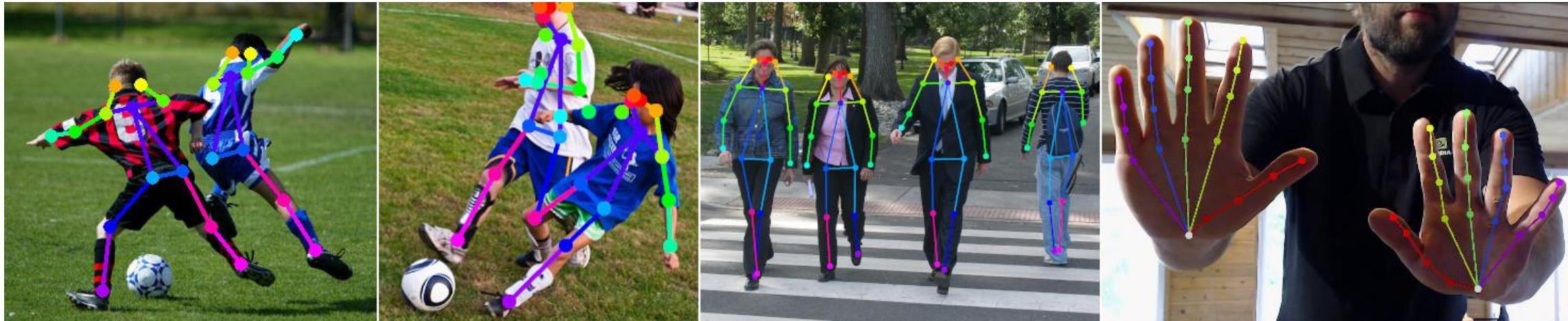


FastDepth, ICRA'19



# Reference – Pose Estimation

Human pose estimation (<https://github.com/dusty-nv/jetson-inference/blob/master/docs/posenet.md>)



→ Locating various body parts (human body + hands) that form a **skeletal topology** (keypoints + links).

- [1] Zhe Cao, et al., "Realtime multi-person 2d pose estimation using part affinity fields." CVPR 2017.
- [2] Bin Xiao , et al., "Simple baselines for human pose estimation and tracking." ECCV 2018.

# Experiments – Depth & Pose Estimation

\*Your basic workspace is here: “`cd ~/jetson-inference/build/aarch64/bin`”

## ### Run mono-depth with image files ###

Q1. Check the original image file “`images/trail_8.jpg`”. Run “`python3 depthnet.py images/trail_8.jpg images/test/output_trail_8.jpg`”. Open the result image file, and what is in the file? Please guess the meaning of the different colors in the output figure.

## ### Live depth estimation ###

Q2. Run “`python3 depthnet.py --network=monodepth-fcn-resnet50 --flip-method=rotate-180`”, Please discuss the quality of the results.

## ### Live pose estimation ###

Q3. Run “`python posenet.py --flip-method=rotate-180 csi://0`”, and “`python posenet.py --network=resnet18-hand --flip-method=rotate-180 csi://0`”. Try different networks and discuss how each output is different. What does each node (keypoint) and link mean?

You can check the models in the below link:

<https://github.com/dusty-nv/jetson-inference/blob/master/docs/posenet.md#pre-trained-pose-estimation-models>

# Experiments – Depth & Pose Estimation

## ### How to try other models? ###

Check your installed model,

```
"cd ~/jetson-inference/build/aarch64/bin/networks"  
"ls"
```

you can see the model folder installed.

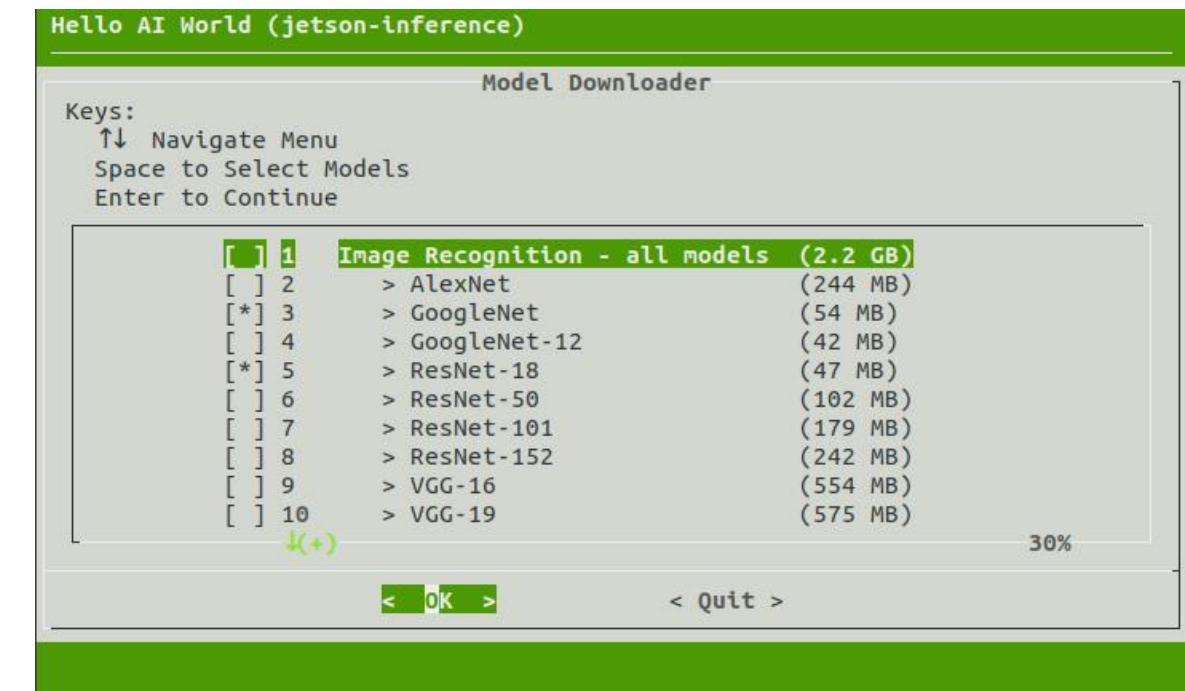
Download models,

```
"cd ~/jetson-inference/tools/"  
"./download-models.sh"
```

After download,

```
"cd ~/jetson-inference/build/aarch64/bin"
```

- Submit all experimental results (classification, detection, segmentation, depth, pose) in **one** report.
- The deadline is postponed until **Nov. 14<sup>th</sup> (Mon)**.



# Tesla's Autopilot is Now on Humanoid Robot

→ Semantic + 3D Visual Perception

In the **Tesla AI Day 2022** (about a month ago),

[https://youtu.be/ODSJsviD\\_SU](https://youtu.be/ODSJsviD_SU)

