



# Robotics: Introduction to perception

Vladimír Petřík

[vladimir.petrik@cvut.cz](mailto:vladimir.petrik@cvut.cz)

23.10.2023

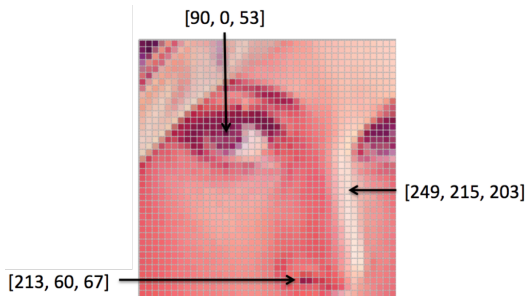
# What is image?

- ▶ Camera connected to computer produces images
- ▶ Image is array of numbers<sup>1</sup>



187	189	174	169	160	162	129	191	172	161	165	166	
185	182	169	16	75	62	39	17	170	219	180	164	
180	180	82	14	54	6	13	33	48	156	159	181	
206	199	6	154	191	111	120	204	146	15	64	180	
164	66	130	281	237	239	239	228	227	67	71	201	
172	156	207	233	233	214	220	239	238	61	74	206	
188	66	178	209	186	216	211	166	191	75	81	168	
189	81	189	14	12	160	126	11	21	62	62	166	
199	146	191	193	198	227	178	185	182	166	96	190	
209	174	199	282	236	231	149	178	238	43	15	204	
196	214	154	149	236	187	66	160	79	36	218	241	
196	224	147	146	227	210	127	161	36	168	285	224	
196	214	173	66	103	143	64	66	7	168	249	218	
187	196	238	76	1	81	47	6	4	237	286	211	
183	202	237	191	6	6	12	108	209	196	143	236	
196	206	123	207	177	121	128	123	203	176	13	64	218

187	189	174	169	160	162	129	191	172	161	165	166	
185	182	169	16	75	62	39	17	170	219	180	164	
180	180	82	14	54	6	13	33	48	156	159	181	
206	199	6	154	191	111	120	204	146	15	64	180	
164	66	130	281	237	239	239	228	227	67	71	201	
172	156	207	233	233	214	220	239	238	61	74	206	
188	66	178	209	186	216	211	166	191	75	81	168	
189	81	189	14	12	160	126	11	21	62	62	166	
199	146	191	193	198	227	178	185	182	166	96	190	
209	174	199	282	236	231	149	178	238	43	15	204	
196	214	154	149	236	187	66	160	79	36	218	241	
196	224	147	146	227	210	127	161	36	168	285	224	
196	214	173	66	103	143	64	66	7	168	249	218	
187	196	238	76	1	81	47	6	4	237	286	211	
183	202	237	191	6	6	12	108	209	196	143	236	
196	206	123	207	177	121	128	123	203	176	13	64	218

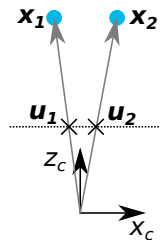
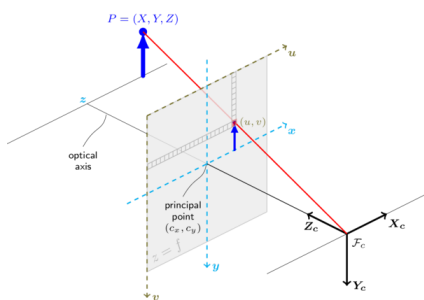


<sup>1</sup>Images are from: <https://ai.stanford.edu/~syueung/cvweb/tutorial1.html>



# How is the image formed?

- ▶ Perspective camera
  - ▶ pinhole camera model<sup>2</sup>
  - ▶ projects spatial point  $x_c$  into image point  $u = (u \ v)^\top$  by intersecting
    - ▶ image plane and
    - ▶ the line connecting  $x_c$  with the projection center
  - ▶ all points on a ray project to the same pixel



<sup>2</sup>[docs.opencv.org](https://docs.opencv.org)



# Projection of pinhole camera

- ▶  $\mathbf{u}_H = K\mathbf{x}_c$ 
  - ▶  $\mathbf{u}_H$  is pixel in homogeneous coordinates
  - ▶ if  $\mathbf{u}_H = (u_H \ v_H \ w_H)^\top$ , then pixel coordinates are  $(u_H/w_H \ v_H/w_H)^\top$
  - ▶ alternatively, we can represent it as:  $\lambda(u, v, 1)^\top = K\mathbf{x}_c$
- ▶  $K$  is camera matrix
  - ▶  $K = \begin{pmatrix} f_x & 0 & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{pmatrix}$
  - ▶ what does  $\lambda$  represent?
    - ▶  $\lambda$  is non-zero real number
    - ▶ if you know  $\lambda$  value, you can compute Cartesian coordinate  $\mathbf{x} = \lambda K^{-1}\mathbf{u}$
    - ▶ otherwise, only ray is computable
  - ▶ how to find  $K$  from points?



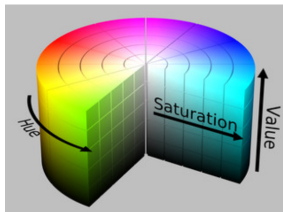
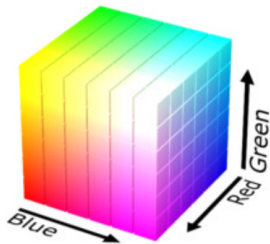
# What we can study on images?

- ▶ Segmentation masks (where are the objects of interest)
- ▶ Objects classification (labeling)



# Segmentation masks - color thresholding

- ▶ Thresholding
  - ▶ RGB pixel values for coordinates  $\mathbf{u}$ :  $I_{\text{RGB}}(\mathbf{u})$
  - ▶  $M(\mathbf{u}) = 1$ , if  $I_{\text{RGB}}(\mathbf{u}) = (0 \ 255 \ 0)^T$  ?
  - ▶  $M(\mathbf{u}) = 1$ , if  $\tau_l < I_{\text{RGB}}(\mathbf{u}) < \tau_u$ , for all channels
  - ▶  $M(\mathbf{u}) = 1$ , if  $\varphi_l < I_{\text{HSV}}(\mathbf{u}) < \varphi_u$ , for all channels
- ▶ Post-processing
  - ▶ compute connected components
  - ▶ remove small or deformed segments
  - ▶ assign label based on thresholds

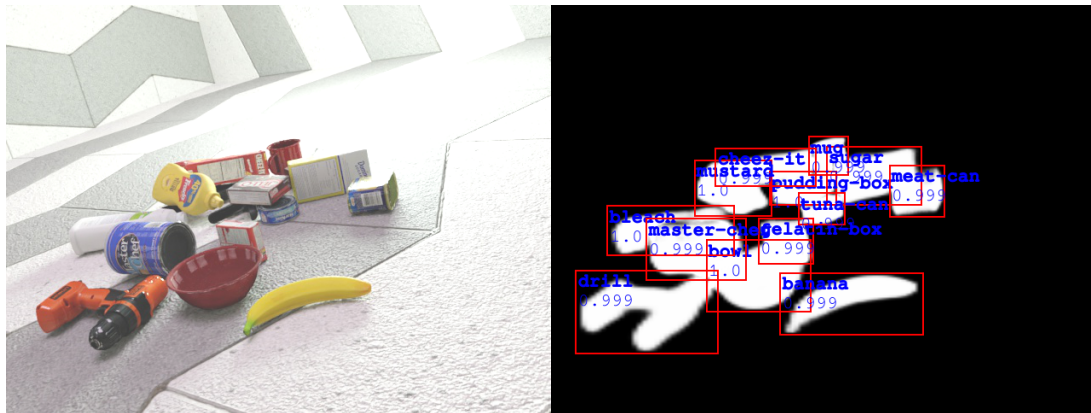


# Segmentation masks for known 3D objects

- ▶ Neural Network (e.g. Mask R-CNN)
- ▶ Training inputs:
  - ▶ dataset of images, masks and labels, or
  - ▶ dataset of known 3D objects (meshes)
  - ▶ quality depends on the training data (augmentations)
- ▶ Inference:
  - ▶ Input: image
  - ▶ Output: segmentation mask, bounding box, label, and confidence

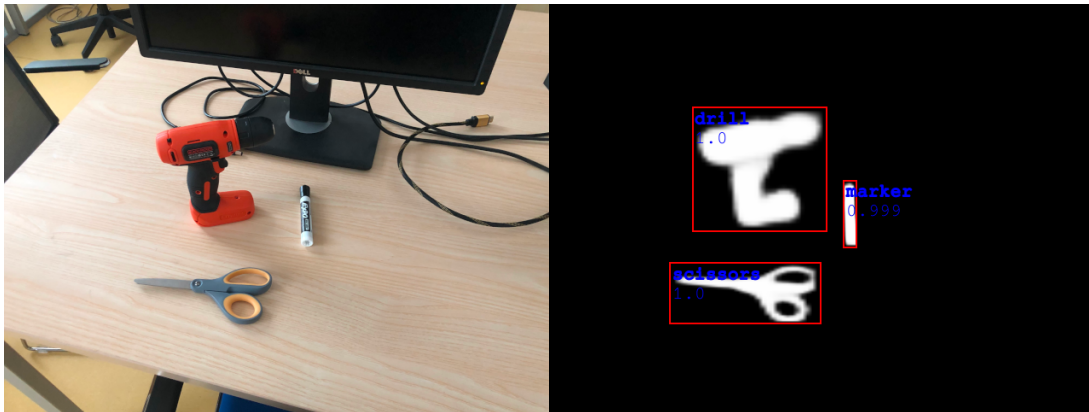


# Mask R-CNN results





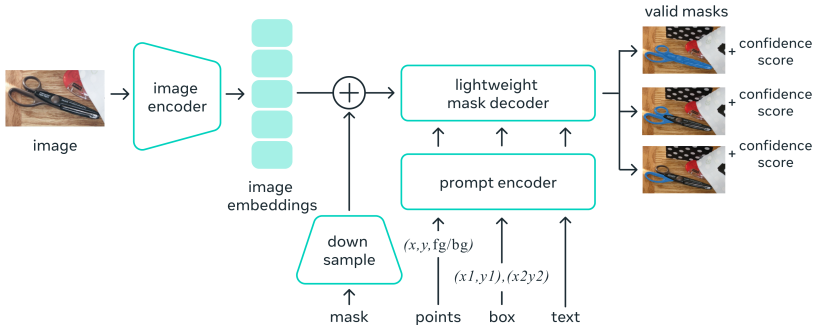
# Mask R-CNN results



# Segmentation masks without re-training

- ▶ Segment Anything Model (SAM)
  - ▶ segment any object, in any image, with a single click
  - ▶ dataset of 10M images, 1B masks

## Universal segmentation model



## SAM results



## SAM results



# Segmentation

- ▶ Segmentation finds objects in image
  - ▶ segmentation mask
  - ▶ bounding box
  - ▶ label
  - ▶ confidence score
- ▶ Information only in image space
- ▶ How to use it in robot space?



## External camera

- ▶ Assume camera mounted rigidly to the reference frame
  - ▶ if we know  $K$  and  $T_{RC}$ , how to project points  $x_R$  to image?
- ▶ Unknown  $K$  and  $T_{RC}$  and planar problem
  - ▶ e.g. cubes with the same high on table desk
  - ▶ what is the position of cube on 2D table w.r.t. 2D image/pixels coordinates?
  - ▶ analyzed by **homography**

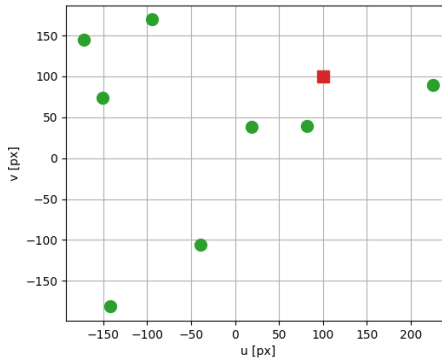
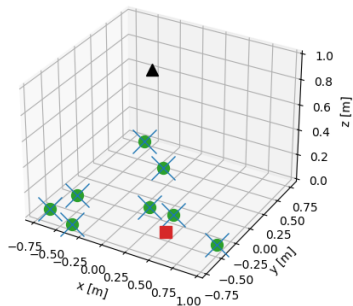


# Homography

- ▶ Homography matrix  $H$  is  $3 \times 3$  matrix that maps points from one plane to another
  - ▶ image plane to table desk
  - ▶ one image plane to another image plane (different view)
- ▶  $s \begin{pmatrix} x & y & 1 \end{pmatrix}^\top = H \begin{pmatrix} u & v & 1 \end{pmatrix}^\top$ 
  - ▶  $x, y$  are coordinates in the first plane
  - ▶  $u, v$  are coordinates in the second plane
- ▶ 9 elements but only 8 DoF, usually added constraint  $h_{33} = 1$
- ▶ How to find H?
  - ▶  $H, _ = \text{cv2.findHomography}(U, X)$
  - ▶  $U, X$  are  $N \times 2$  correspondence points
  - ▶ e.g. measure manually
    - ▶ position of cube center w.r.t. table corner
    - ▶ position of cube center in image



# Homography example





# Non-planar pose estimation

- ▶ Homography maps only plane to plane
- ▶ More general object pose estimation in **camera** frame
  - ▶ get depth by mapping from area in pixels to depth for fixed size objects
  - ▶ get depth by additional scene information, e.g. known size/model of the objects
  - ▶ RGBD camera
  - ▶ additional markers





# Using depth sensor

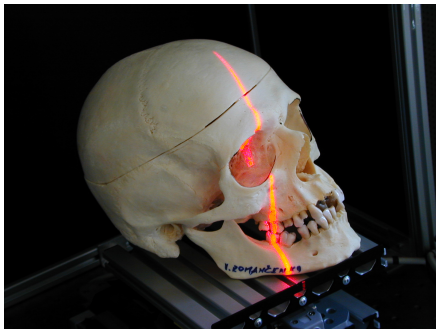
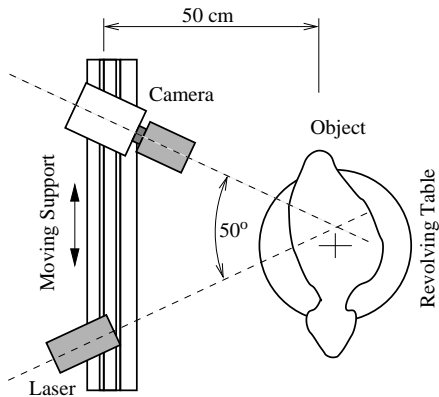
## ▶ RGBD sensors

- ▶ RGB image ( $H \times W \times 3$ )
- ▶ Depth map ( $H \times W \times 1$ ), distance in meters for each pixel
- ▶ Structured point cloud ( $H \times W \times 3$ ),  $(x_c \ y_c \ z_c)$  for each pixel



## How depth sensor works

- ▶ Laser projects pattern and camera recognizes it
- ▶ Depth information is computed using triangulation



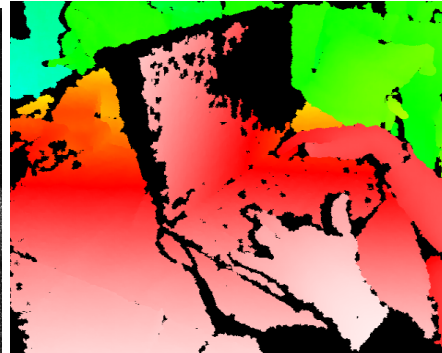
## 2D depth sensors

- ▶ Based on the structured light
- ▶ Projects 2D infra red patterns
- ▶ One projector and two cameras (RGB + IR)



## Issues with depth sensors

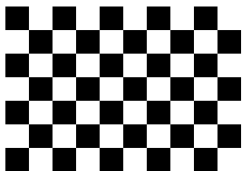
- ▶ Depth reconstruction is not perfect (black areas in the image<sup>3</sup>)
- ▶ In python represented by NaN
- ▶ Not every pixel in RGB has reconstructed depth value
- ▶ RGB and Depth data are not aligned (you need to calibrate them)



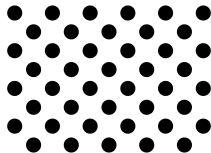
<sup>3</sup><https://commons.wikimedia.org>, User:Kolossos

## Additional markers

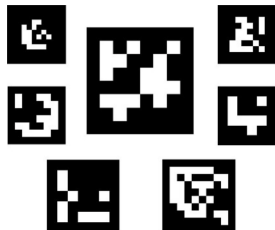
- ▶ Can we compute the pose of patterns<sup>4</sup>?
  - ▶ the size and structure needs to be known
  - ▶ subpixel accuracy
  - ▶ it has to be completely visible
- ▶ Can we compute the pose of ArUco markers?
  - ▶ less accurate than regular patterns
  - ▶ provides marker id and the pose
  - ▶ it has to be completely visible



© 2008 Intel Corporation  
All rights reserved.



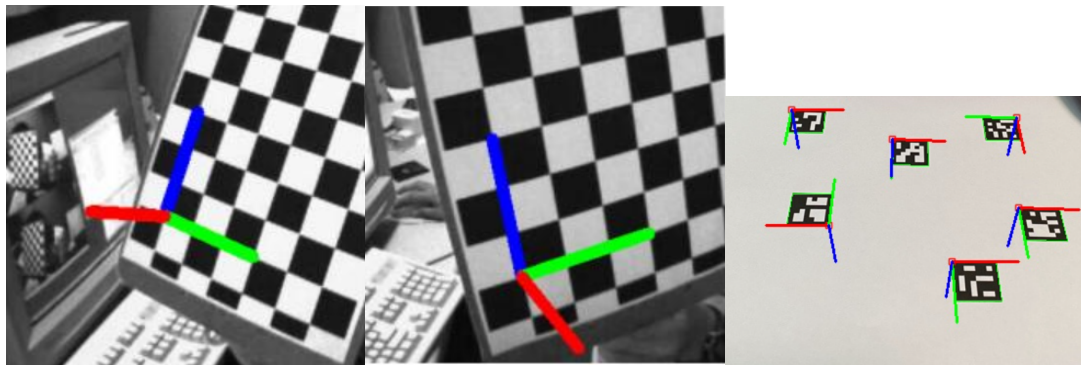
© 2008 Intel Corporation  
All rights reserved.



<sup>4</sup>[docs.opencv.org](http://docs.opencv.org)



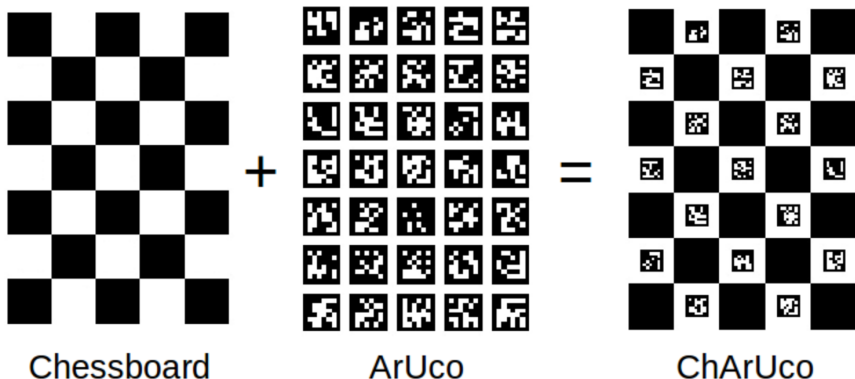
## Markers pose example





## ChArUco board for calibration

- ▶ Combines accuracy of pattern with detections of ArUco
- ▶ Partial visibility detections



Charuco definition



## Camera matrix estimation with boards

- ▶ We can estimate camera matrix from correspondences in image space and spatial space
  - ▶ collect images of the board from different views
  - ▶ detect boards
  - ▶ compute correspondences between image points and board frame points
  - ▶ `_, K, dist_coeffs, rvecs, tvecs = cv2.calibrateCamera(obj_points, img_points, img_shape)`
- ▶ In addition we get
  - ▶ distortion coefficients that compensates defects of objective
    - `Knew, roi = cv.getOptimalNewCameraMatrix(K, dist_coeffs, img_shape, 1, img_shape)`
    - `img_undistorted = cv.undistort(img, K, dist_coeffs, None, Knew)`
  - ▶  $SE(3)$  poses of boards in camera frame



# Pose estimation from RGB(D)

- ▶ Pose estimation methods
  - ▶ use prior knowledge about the task, e.g. fixed height objects on a plane
  - ▶ use prior knowledge about the objects (size)
  - ▶ use depth sensor
  - ▶ use ArUco markers
- ▶ Where is robot?
  - ▶ homography estimates poses of objects w.r.t. plane frame
  - ▶ other methods estimate poses in camera frame
  - ▶ we need to estimate/calibrate  $T_{RC}$



# HandEye calibration

- ▶ Camera can be mounted w.r.t.
  - ▶ robot base frame (eye-to-hand calibration)
  - ▶ gripper frame (eye-in-hand calibration)
- ▶ Solve  $A^i X = Y B^i$ 
  - ▶ measurements:  $A^i, B^i \in SE(3)$
  - ▶ estimated parameters:  $X, Y \in SE(3)$
- ▶  $X, Y = \text{calibrateRobotWorldHandEye}(A, B)$
- ▶ Eye-to-hand calibration
  - ▶  $A^i = T_{RG}^i$
  - ▶  $B^i = T_{CT}^i$
  - ▶  $X = T_{GT}$
  - ▶  $Y = T_{RC}$
- ▶ Eye-in-hand calibration
  - ▶  $A^i = T_{CT}^i$
  - ▶  $B^i = T_{GR}^i$
  - ▶  $X = T_{TR}$
  - ▶  $Y = T_{CG}$



# Summary

- ▶ Image representation
- ▶ Projection to/from image
- ▶ Segmentation in image space
- ▶ Homography
- ▶ Pose estimation from image
- ▶ Camera calibration



# Laboratory

- ▶ No new homework this week
- ▶ Homography estimation on toy example in Python/OpenCV
- ▶ HandEye calibration on toy example in Python/OpenCV

