



# Robotics: Introduction to perception

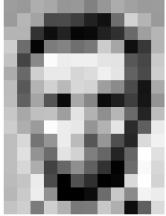
Vladimír Petřík

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

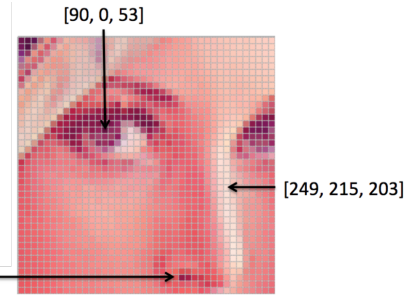
23.10.2023

# What is image?

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



187	180	174	168	162	156	150	144	138	132	126	120	114	108	102
185	182	168	154	144	134	124	114	104	94	84	74	64	54	44
180	180	80	14	14	4	4	30	30	40	50	60	70	80	90
204	174	4	100	110	120	130	140	150	160	170	180	190	200	210
164	80	100	100	100	100	100	100	100	100	100	100	100	100	100
172	100	200	200	210	220	230	240	250	260	270	280	290	300	310
188	88	178	208	182	214	211	184	198	176	20	198	176	20	198
188	97	168	84	168	168	154	11	21	62	22	168	168	168	168
199	186	191	199	188	227	178	181	182	181	181	181	181	181	181
209	178	191	232	204	201	184	178	228	40	112	204	184	184	184
166	216	182	168	168	168	168	168	168	168	168	168	168	168	168
166	214	147	166	227	212	177	161	161	161	161	161	161	161	161
166	214	173	84	162	142	14	36	1	100	168	214	168	214	168
187	186	208	16	1	81	47	9	1	217	208	211	208	211	208
183	202	207	161	8	6	14	188	204	204	204	204	204	204	204
186	206	202	207	177	185	188	208	176	183	188	216	186	216	186

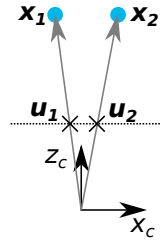
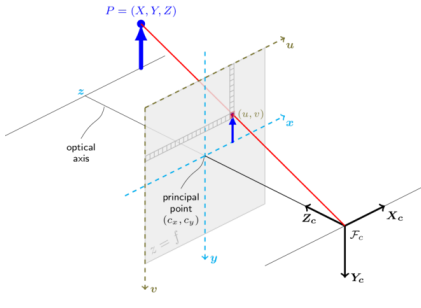


<sup>1</sup>Images are from: <https://ai.stanford.edu/~syyeung/cvweb/tutorial1.html>



# 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



# Projection of pinhole camera

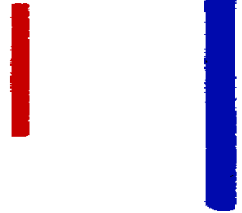
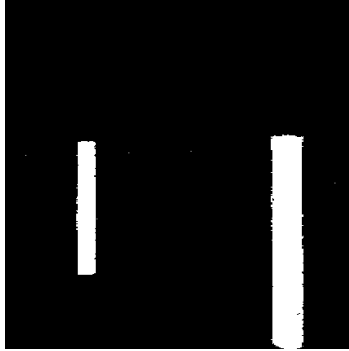
- ▶  $\mathbf{u}_H = K\mathbf{x}_c$ 
  - ▶  $\mathbf{u}_H$  is pixel in homogeneous coordinates
  - ▶ if  $\mathbf{u}_H = (u_H \ v_H \ w_H)^\top$ , then pixel coordinates are  $(u_H/w_H \ v_H/w_H)^\top$
  - ▶ alternatively, we can represent it as:  $\lambda(u, v, 1)^\top = K\mathbf{x}_c$
- ▶  $K$  is 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?





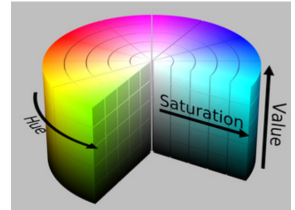
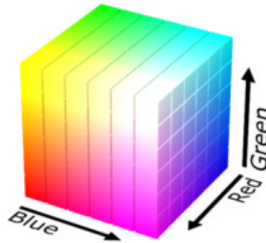
## What we can study on images?

- ▶ 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

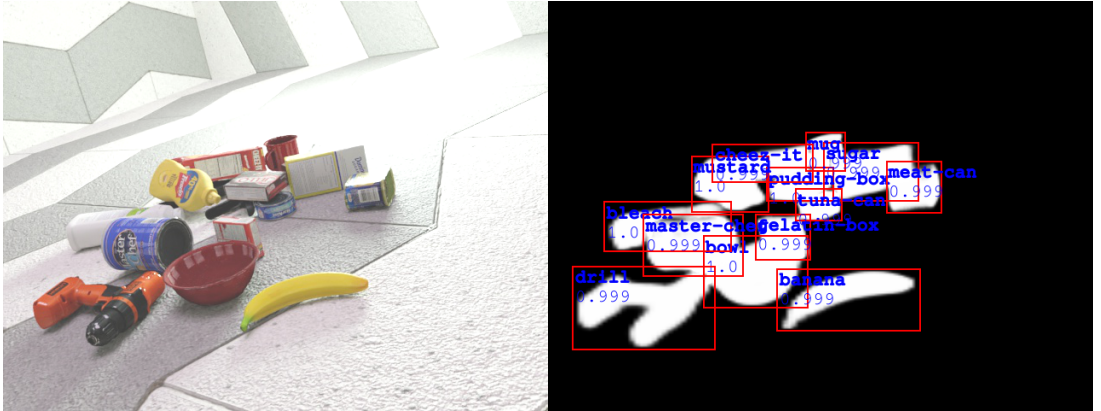


# Segmentation masks for known 3D objects

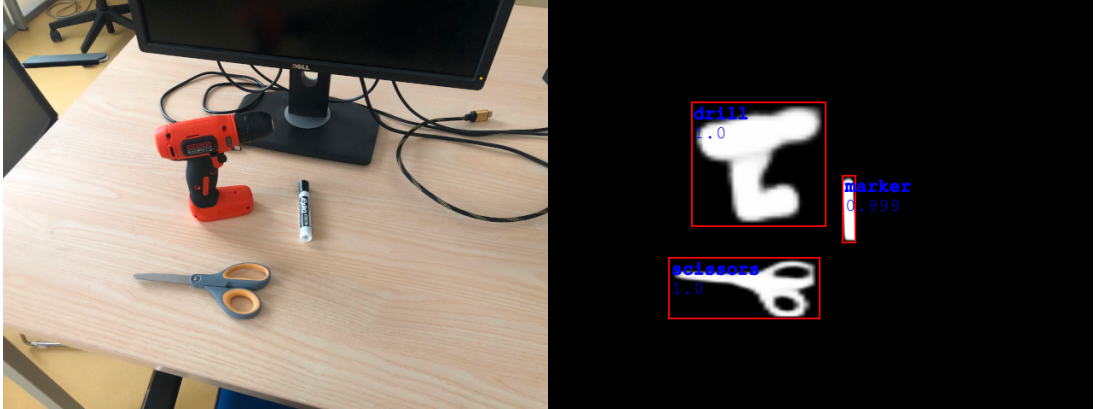
- ▶ Neural Network (e.g. Mask R-CNN)
- ▶ Training inputs:
  - ▶ dataset of images, masks and labels, or
  - ▶ dataset of known 3D objects (meshes)
  - ▶ quality depends on the training data (augmentations)
- ▶ Inference:
  - ▶ Input: image
  - ▶ Output: segmentation mask, bounding box, label, and confidence



# Mask R-CNN results



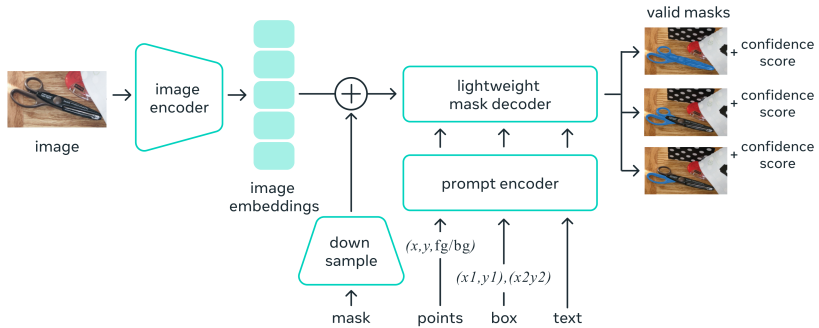
# Mask R-CNN results



# 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

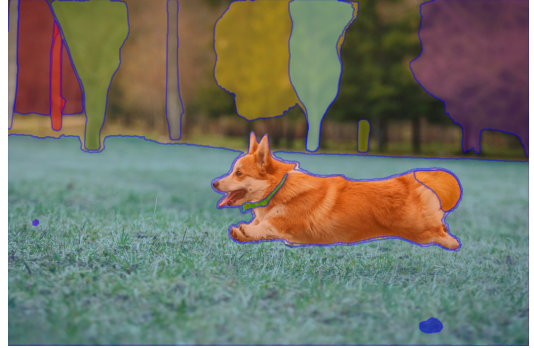
## Universal segmentation model



# 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**

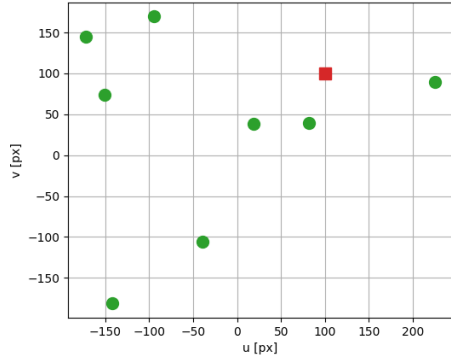
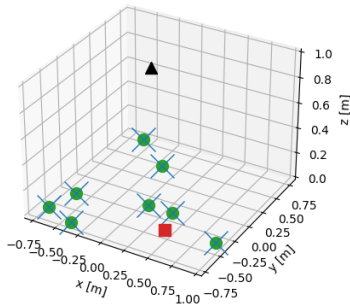


# Homography

- ▶ 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



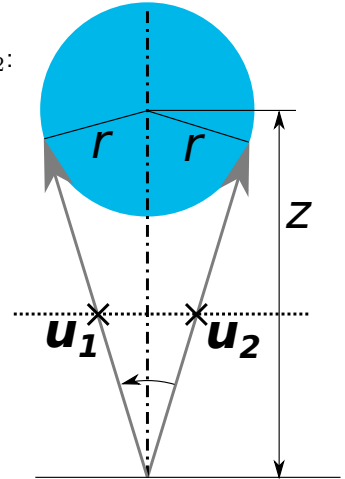
# Non-planar pose estimation

- ▶ 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 prior knowledge about size

- ▶ We know radius is fixed
- ▶ From detected pixels  $\mathbf{u}_1, \mathbf{u}_2$ , we can compute rays  $\mathbf{x}_1, \mathbf{x}_2$ :  
$$\frac{1}{\lambda_i} \mathbf{x}_i = K^{-1} \mathbf{u}_i$$
- ▶ Angle between vectors:  $\cos \alpha = \frac{\frac{1}{\lambda_1 \lambda_2} \mathbf{x}_1 \cdot \mathbf{x}_2}{\frac{1}{\lambda_1 \lambda_2} \|\mathbf{x}_1\| \|\mathbf{x}_2\|}$
- ▶ Depth:  $z = \frac{r}{\sin(\alpha/2)}$



# Using depth sensor

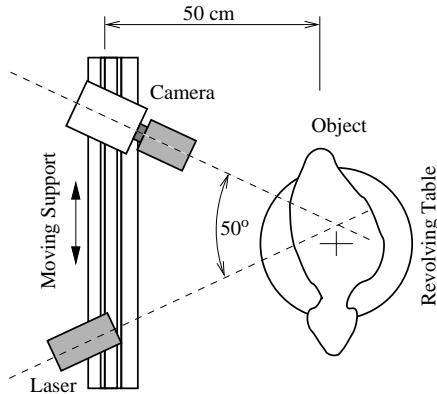
## ▶ RGBD sensors

- ▶ 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





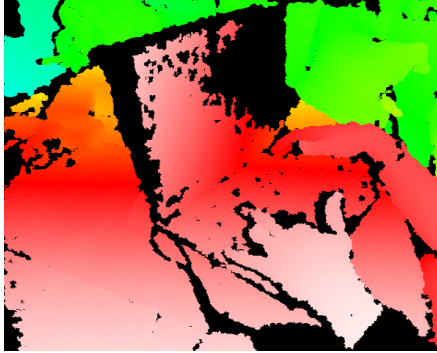
## 2D depth sensors

- ▶ 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)

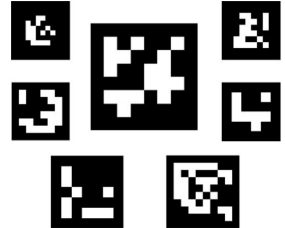
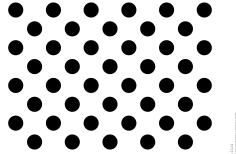
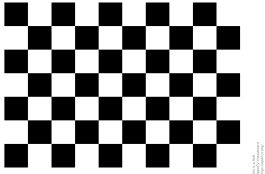


<sup>3</sup><https://commons.wikimedia.org>, User:Kolossos



## Additional markers

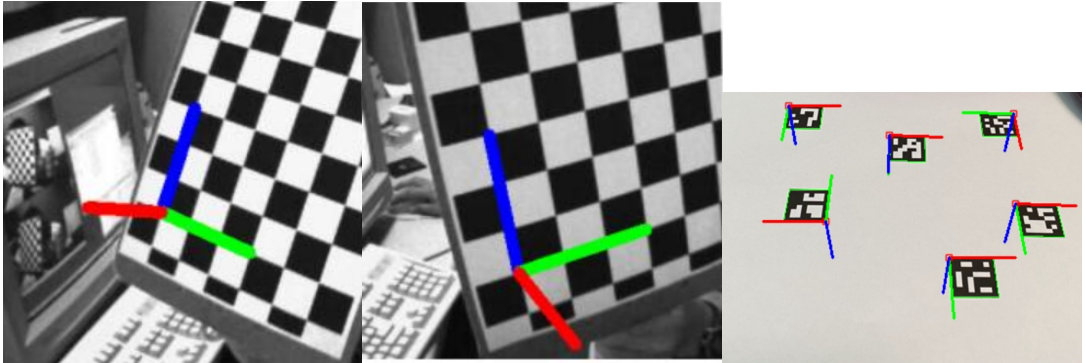
- ▶ Can we compute the pose of patterns<sup>4</sup>?
  - ▶ the size and structure needs to be known
  - ▶ subpixel accuracy
  - ▶ it has to be completely visible
- ▶ Can we compute the pose of ArUco markers?
  - ▶ less accurate than regular patterns
  - ▶ provides marker id and the pose
  - ▶ it has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



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

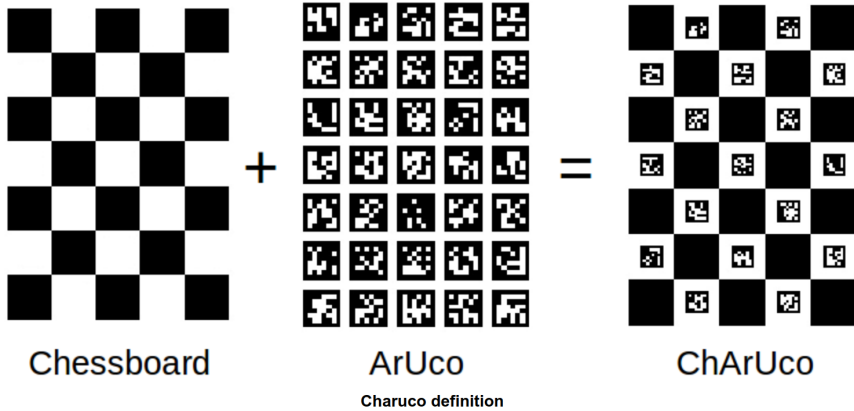


# Markers pose example



# ChArUco board for calibration

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



# 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

