



Robotics: Introduction to perception

Vladimír Petřík

vladimir.petrík@cvut.cz

21.10.2023

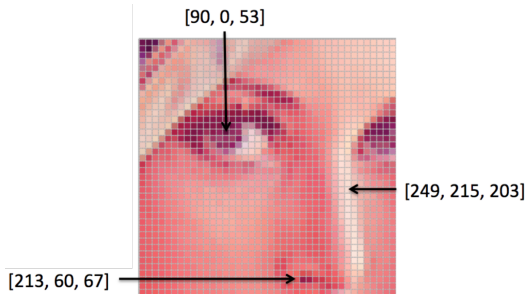
What is image?

- ▶ Camera connected to computer produces images
- ▶ Image is array of numbers¹



187	189	174	148	160	162	120	101	170	163	165	166
185	182	169	14	75	63	55	17	170	210	180	164
180	180	82	14	34	6	10	33	40	156	169	181
206	100	6	154	131	111	120	254	146	16	64	180
194	68	137	281	237	239	239	228	227	87	71	201
172	106	207	233	233	214	220	239	238	58	74	206
188	88	178	209	185	216	211	158	189	75	50	188
189	97	168	14	12	160	126	11	21	62	22	148
199	148	191	193	198	227	178	143	182	106	86	190
205	174	191	292	236	231	149	178	238	43	53	204
190	214	116	149	236	187	86	160	79	38	218	241
190	224	147	186	227	210	127	103	36	108	285	224
190	214	173	86	103	143	96	50	2	108	249	215
187	196	238	76	1	81	47	5	6	217	286	211
183	202	237	195	6	0	12	108	200	138	243	236
196	206	123	207	177	121	139	200	178	18	81	218

187	189	174	148	160	162	120	101	170	163	165	166
185	182	169	14	75	63	55	17	170	210	180	164
180	180	82	14	34	6	10	33	40	156	169	181
206	100	6	154	131	111	120	254	146	16	64	180
194	68	137	281	237	239	239	228	227	87	71	201
172	106	207	233	233	214	220	239	238	58	74	206
188	88	178	209	185	216	211	158	189	75	50	188
189	97	168	14	12	160	126	11	21	62	22	148
199	148	191	193	198	227	178	143	182	106	86	190
205	174	191	292	236	231	149	178	238	43	53	204
190	214	116	149	236	187	86	160	79	38	218	241
190	224	147	186	227	210	127	103	36	108	285	224
190	214	173	86	103	143	96	50	2	108	249	215
187	196	238	76	1	81	47	5	6	217	286	211
183	202	237	195	6	0	12	108	200	138	243	236
196	206	123	207	177	121	139	200	178	18	81	218

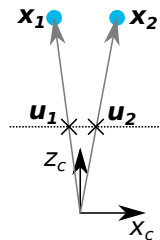
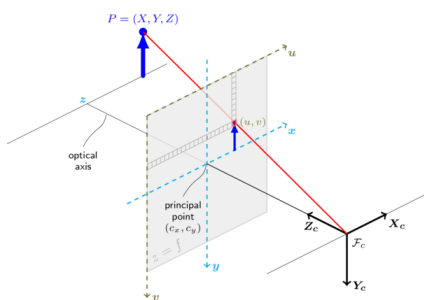


¹Images are from: <https://ai.stanford.edu/~syueung/cvweb/tutorial1.html>



How is the image formed?

- ▶ Perspective camera
 - ▶ pinhole camera model²
 - ▶ projects spatial point x_c into image point $u = \begin{pmatrix} u & v \end{pmatrix}^\top$ by intersecting
 - ▶ image plane and
 - ▶ the line connecting x_c with the projection center
 - ▶ all points on a ray project to the same pixel



²docs.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 λ represent?
 - ▶ λ is non-zero real number
 - ▶ if you know λ value, you can compute Cartesian coordinate $\mathbf{x} = \lambda K^{-1} \mathbf{u}$
 - ▶ otherwise, only ray is computable
 - ▶ how to find K from points?



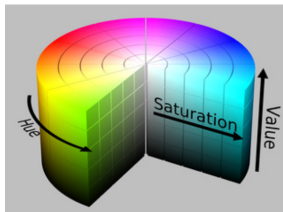
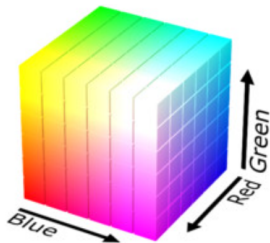
What we can study on images?

- ▶ 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 $I_{\text{RGB}}(\mathbf{u}) = (0 \ 255 \ 0)^{\top}$?
 - ▶ $M(\mathbf{u}) = 1$, 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

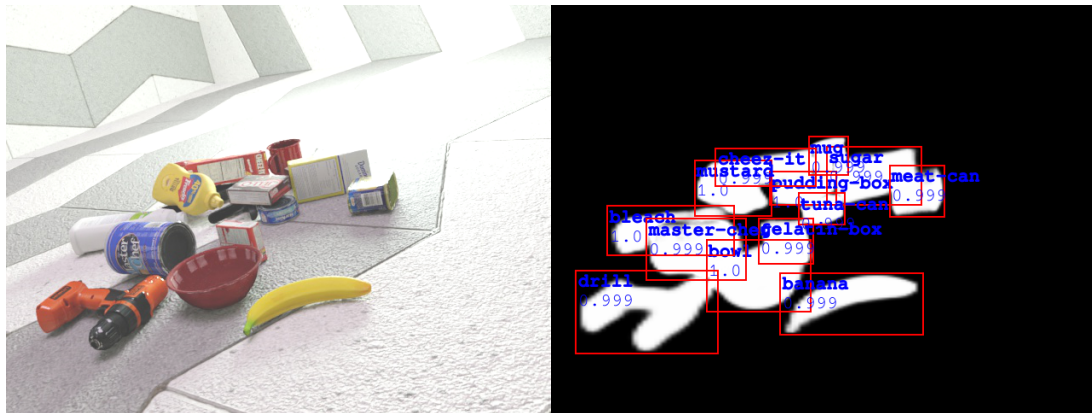


Segmentation masks for known 3D objects

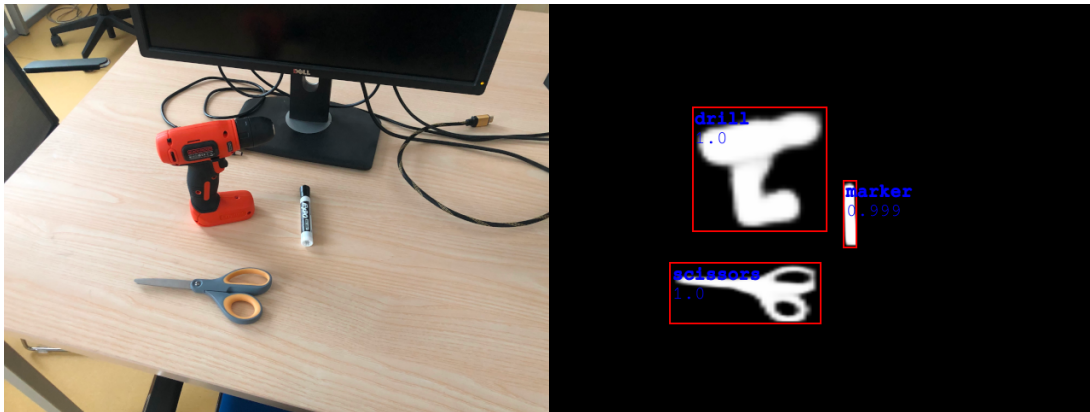
- ▶ Neural Network (e.g. Mask R-CNN)
- ▶ Training inputs:
 - ▶ dataset of images, masks and labels, or
 - ▶ dataset of known 3D objects (meshes)
 - ▶ quality depends on the training data (augmentations)
- ▶ Inference:
 - ▶ Input: image
 - ▶ Output: segmentation mask, bounding box, label, and confidence



Mask R-CNN results



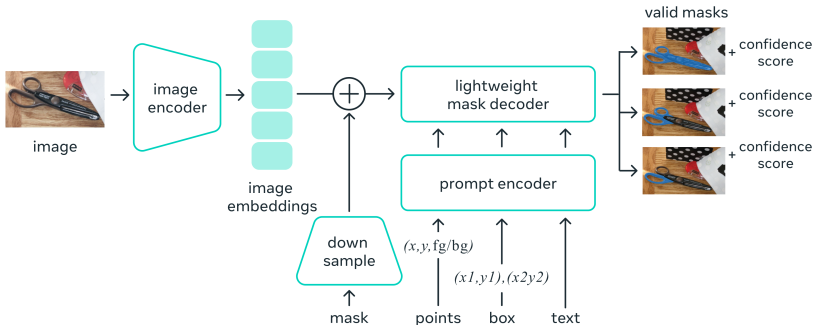
Mask R-CNN results



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

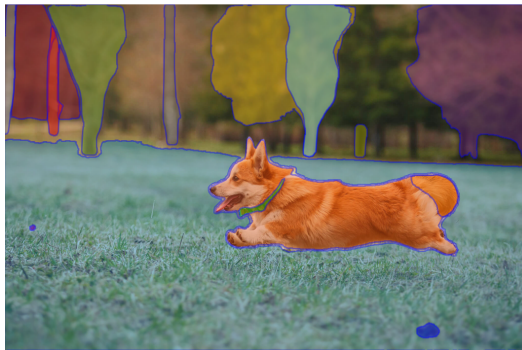
Universal segmentation model

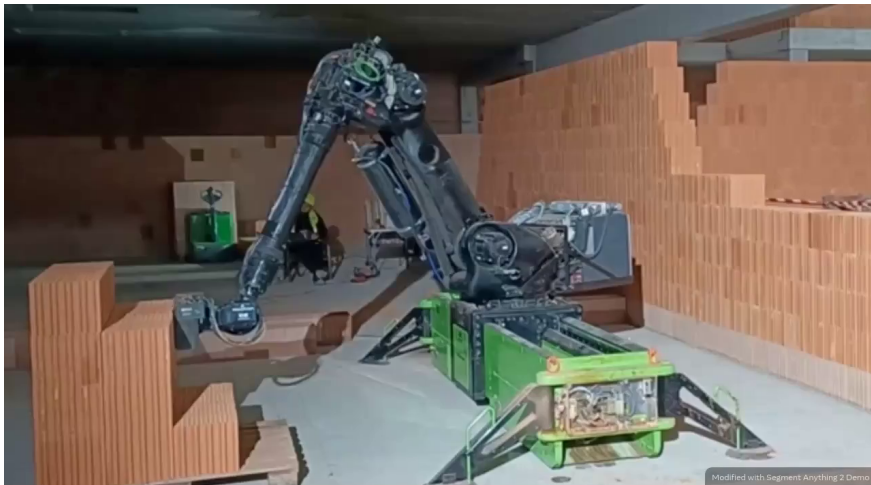


SAM results



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**

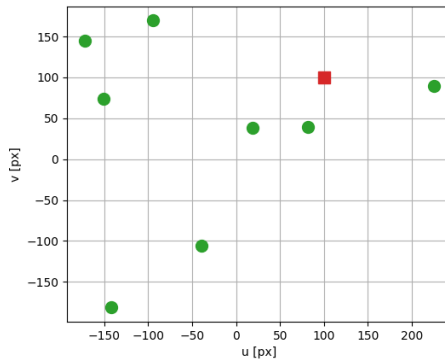
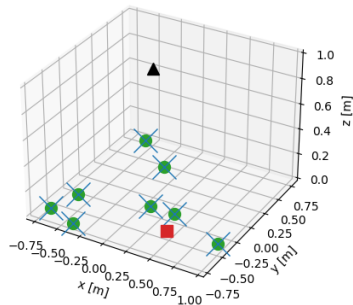


Homography

- ▶ Homography matrix H is 3×3 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



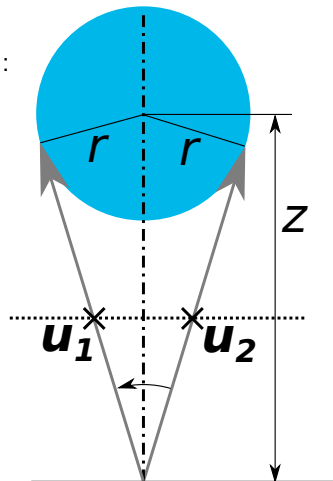
Non-planar pose estimation

- ▶ 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: $\cos \alpha = \frac{\frac{1}{\lambda_1 \lambda_2} \mathbf{x}_1 \cdot \mathbf{x}_2}{\frac{1}{\lambda_1} \|\mathbf{x}_1\| \frac{1}{\lambda_2} \|\mathbf{x}_2\|}$
- ▶ Depth: $z = \frac{r}{\sin(\alpha/2)}$



Using depth sensor

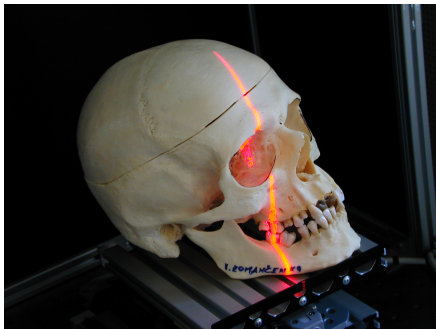
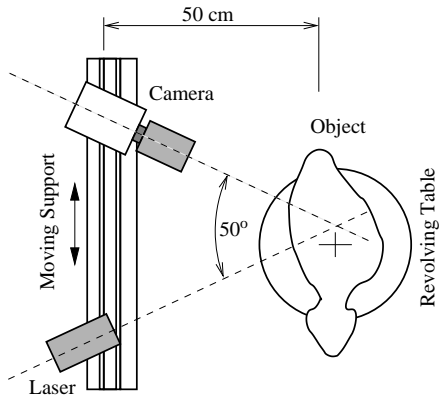
► RGBD sensors

- 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



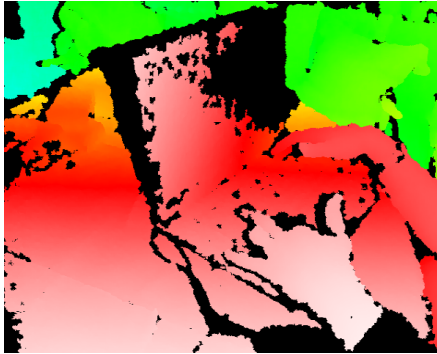
2D depth sensors

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



Issues with depth sensors

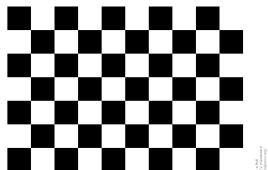
- ▶ Depth reconstruction is not perfect (black areas in the image³)
- ▶ In python represented by NaN
- ▶ Not every pixel in RGB has reconstructed depth value
- ▶ RGB and Depth data are not aligned (you need to calibrate them)



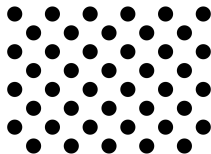
³<https://commons.wikimedia.org>, User:Kolossos

Additional markers

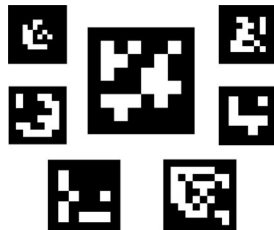
- ▶ Can we compute the pose of patterns⁴?
 - ▶ the size and structure needs to be known
 - ▶ subpixel accuracy
 - ▶ it has to be completely visible
- ▶ Can we compute the pose of ArUco markers?
 - ▶ less accurate than regular patterns
 - ▶ provides marker id and the pose
 - ▶ it has to be completely visible



OpenCV

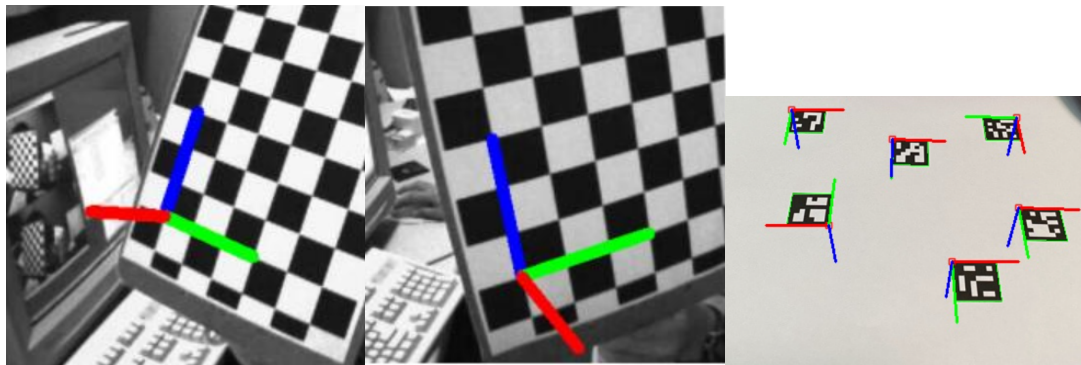


OpenCV



⁴docs.opencv.org

Markers pose example

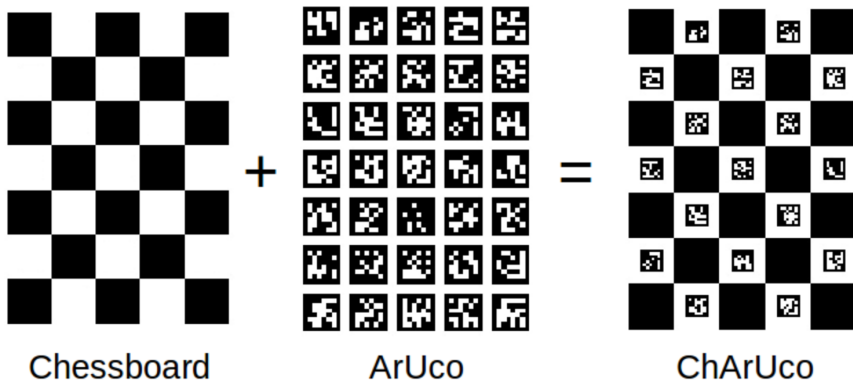


Markers example from real-world



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}



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 in OpenCV
- ▶ Camera calibration in OpenCV



What is next?

- ▶ No lecture/exercise next week
- ▶ In two weeks, we will start with test
 - ▶ Questions from first three lectures
 - ▶ SE2 and SE3 transformations and properties
 - ▶ Forward kinematics and DH notation
 - ▶ Jacobian, its properties, and usage

