

Robotics: Introduction to perception

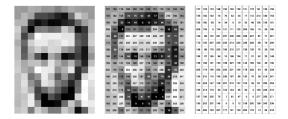
Vladimír Petrík

vladimir.petrik@cvut.cz

6.10.2025

What is image?

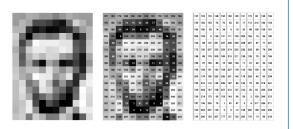
- Camera connected to computer produces images
- ► Image is array of numbers¹

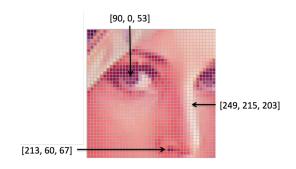


¹Images are from: https://ai.stanford.edu/~syyeung/cvweb/tutorial1.html

What is image?

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

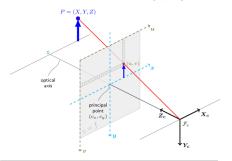




¹Images are from: https://ai.stanford.edu/~syyeung/cvweb/tutorial1.html

How is the image formed?

- Perspective camera
 - pinhole camera model²
 - $lackbox{ projects spatial point } oldsymbol{x}_c \text{ into image point } oldsymbol{u} = egin{pmatrix} u & v \end{pmatrix}^{ op} \text{ by intersecting}$
 - image plane and
 - ightharpoonup the line connecting $oldsymbol{x}_c$ with the projection center
 - all points on a ray project to the same pixel

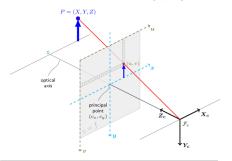


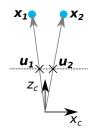
²docs.opencv.org



How is the image formed?

- Perspective camera
 - pinhole camera model²
 - $lackbox{ projects spatial point } oldsymbol{x}_c \text{ into image point } oldsymbol{u} = egin{pmatrix} u & v \end{pmatrix}^{ op} \text{ by intersecting}$
 - image plane and
 - ightharpoonup the line connecting $oldsymbol{x}_c$ with the projection center
 - ▶ all points on a ray project to the same pixel





²docs.opencv.org



- $\mathbf{u}_H = K \mathbf{x}_c$
 - $ightharpoonup u_H$ is pixel in homogeneous coordinates
 - $lackbox{ if } oldsymbol{u}_H = egin{pmatrix} u_H & v_H & v_H \end{pmatrix}^{ op}, \text{ then pixel coordinates are } egin{pmatrix} u_H/w_H & v_H/w_H \end{pmatrix}^{ op}$

- $\mathbf{u}_H = K \mathbf{x}_c$
 - $lackbox{m{u}}_H$ is pixel in homogeneous coordinates
 - ▶ if $u_H = \begin{pmatrix} u_H & v_H & w_H \end{pmatrix}^{\top}$, then pixel coordinates are $\begin{pmatrix} u_H/w_H & v_H/w_H \end{pmatrix}^{\top}$
 - ightharpoonup alternatively, we can represent it as: $\lambda \left(u,v,1\right)^{\top}=K\boldsymbol{x}_{c}$

- $\mathbf{u}_H = K \mathbf{x}_c$
 - $lackbox{ } u_H$ is pixel in homogeneous coordinates
 - ▶ if $u_H = \begin{pmatrix} u_H & v_H & w_H \end{pmatrix}^{\top}$, then pixel coordinates are $\begin{pmatrix} u_H/w_H & v_H/w_H \end{pmatrix}^{\top}$
 - lacktriangle alternatively, we can represent it as: $\lambda \left(u,v,1\right)^{ op}=Koldsymbol{x}_c$
- ightharpoonup K is camera matrix

- $\mathbf{u}_H = K \mathbf{x}_c$
 - $lackbox{m{u}}_H$ is pixel in homogeneous coordinates
 - $lackbox{ if } oldsymbol{u}_H = egin{pmatrix} u_H & v_H & w_H \end{pmatrix}^ op$, then pixel coordinates are $egin{pmatrix} u_H/w_H & v_H/w_H \end{pmatrix}^ op$
 - lacktriangle alternatively, we can represent it as: $\lambda \left(u,v,1\right)^{ op}=Koldsymbol{x}_c$
- K is camera matrix

 \blacktriangleright what does λ represent?

- $\mathbf{u}_H = K \mathbf{x}_c$
 - $lackbox{m{u}}_H$ is pixel in homogeneous coordinates
 - $lackbox{ if } oldsymbol{u}_H = egin{pmatrix} u_H & v_H & v_H \end{pmatrix}^{ op}, \text{ then pixel coordinates are } egin{pmatrix} u_H/w_H & v_H/w_H \end{pmatrix}^{ op}$
 - lacktriangle alternatively, we can represent it as: $\lambda \left(u,v,1\right)^{\top} = K oldsymbol{x}_c$
- ► *K* is camera matrix

$$K = \begin{pmatrix} f_x & 0 & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{pmatrix}$$

- \blacktriangleright what does λ represent?
 - $\triangleright \lambda$ is non-zero real number
 - if you know λ value, you can compute Cartesian coordinate $x = \lambda K^{-1}u$
 - otherwise, only ray is computable

- $\mathbf{u}_H = K \mathbf{x}_c$
 - $lackbox{m u}_H$ is pixel in homogeneous coordinates
 - $lackbox{ if } oldsymbol{u}_H = egin{pmatrix} u_H & v_H & w_H \end{pmatrix}^ op$, then pixel coordinates are $egin{pmatrix} u_H/w_H & v_H/w_H \end{pmatrix}^ op$
 - lacktriangle alternatively, we can represent it as: $\lambda \left(u,v,1\right)^{\top} = K oldsymbol{x}_c$
- ► *K* is camera matrix

$$K = \begin{pmatrix} f_x & 0 & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{pmatrix}$$

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

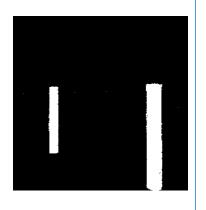
What we can study on images?



What we can study on images?

Segmentation masks (where are the objects of interest)

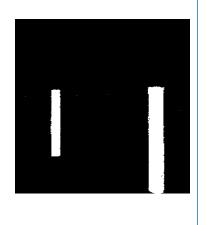




What we can study on images?

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







- ► Thresholding
 - ▶ RGB pixel values for coordinates u: $I_{RGB}(u)$

- ► Thresholding

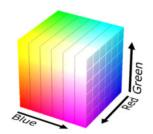
 - ▶ RGB pixel values for coordinates u: $I_{\text{RGB}}(u)$ ▶ M(u) = 1, if $I_{\text{RGB}}(u) = \begin{pmatrix} 0 & 255 & 0 \end{pmatrix}^{\top}$?

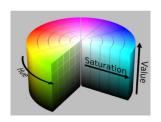
- ► Thresholding
 - ▶ RGB pixel values for coordinates u: $I_{RGB}(u)$
 - M(u) = 1, if $I_{RGB}(u) = \begin{pmatrix} 0 & 255 & 0 \end{pmatrix}^{\top}$?
 - M(u) = 1, if $\tau_l < I_{\mathsf{RGB}}(u) < \tau_u$, for all channels



- Thresholding
 - ▶ RGB pixel values for coordinates u: $I_{RGB}(u)$
 - $M(\boldsymbol{u}) = 1$, if $I_{\mathsf{RGB}}(\boldsymbol{u}) = \begin{pmatrix} 0 & 255 & 0 \end{pmatrix}^{\top}$?
 - M(u) = 1, if $\tau_l < I_{\mathsf{RGB}}(u) < \tau_u$, for all channels
 - $lackbox{M}(oldsymbol{u}) = 1, \text{ if } oldsymbol{arphi}_l < I_{\mathsf{HSV}}(oldsymbol{u}) < oldsymbol{arphi}_u, \text{ for all channels}$

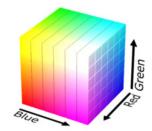


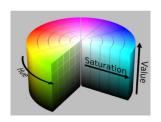




- Thresholding
 - ▶ RGB pixel values for coordinates u: $I_{RGB}(u)$
 - $M(\boldsymbol{u}) = 1$, if $I_{\mathsf{RGB}}(\boldsymbol{u}) = \begin{pmatrix} 0 & 255 & 0 \end{pmatrix}^{\top}$?
 - M(u) = 1, if $\tau_l < I_{\mathsf{RGB}}(u) < \tau_u$, for all channels
- Post-processing
 - compute connected components
 - remove small or deformed segments
 - assign label based on thresholds







► Neural Network (e.g. Mask R-CNN)

- ► Neural Network (e.g. Mask R-CNN)
- ► Training inputs:
 - dataset of images, masks and labels, or

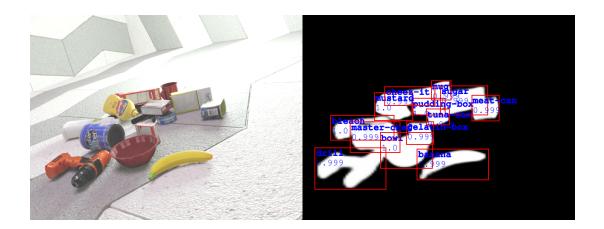
- ► Neural Network (*e.g.* Mask R-CNN)
- ► Training inputs:
 - dataset of images, masks and labels, or
 - dataset of known 3D objects (meshes)

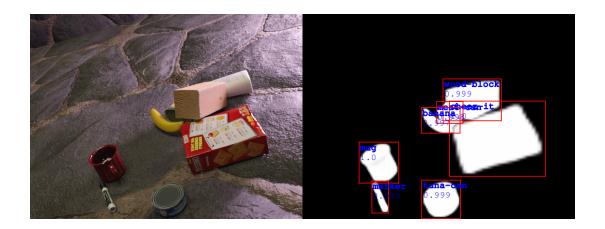


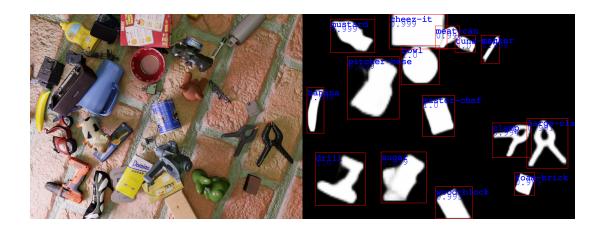


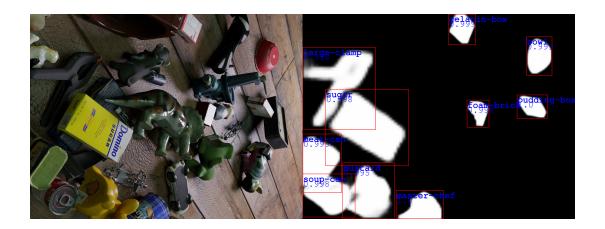
- ► Neural Network (e.g. Mask R-CNN)
- ► Training inputs:
 - dataset of images, masks and labels, or
 - dataset of known 3D objects (meshes)
 - quality depends on the training data (augumentations)

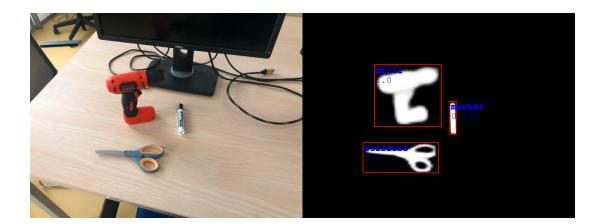
- Neural Network (e.g. Mask R-CNN)
- ► Training inputs:
 - dataset of images, masks and labels, or
 - dataset of known 3D objects (meshes)
 - quality depends on the training data (augumentations)
- Inference:
 - Input: image
 - Output: segmentation mask, bounding box, label, and confidence









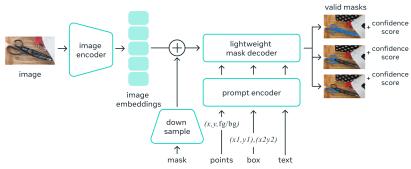




Segmentation masks without re-training

- 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





SAM results



SAM results



SAMv2



Segmentation

- Segmentation finds objects in image
 - segmentation mask
 - bounding box
 - label
 - confidence score

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?

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?

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

- lacktriangle 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)

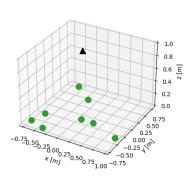
- lacktriangle Homography matrix H is 3×3 matrix that maps points from one plane to another
 - image plane to table desk
 - one image plane to another image plane (different view)
- - ightharpoonup x, y are coordinates in the first plane
 - ightharpoonup u, v are coordinates in the second plane

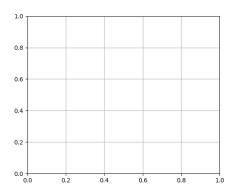
- lacktriangle Homography matrix H is 3×3 matrix that maps points from one plane to another
 - image plane to table desk
 - one image plane to another image plane (different view)
- - ightharpoonup x, y are coordinates in the first plane
 - ightharpoonup u,v are coordinates in the second plane
- ▶ 9 elements but only 8 DoF, usually added constraint $h_{33} = 1$
- ► How to find H?

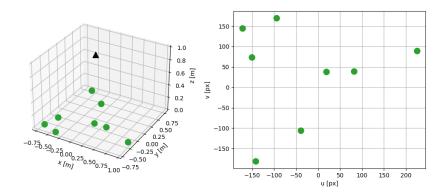
- lacktriangle Homography matrix H is 3×3 matrix that maps points from one plane to another
 - image plane to table desk
 - one image plane to another image plane (different view)
- $ightharpoonup s (x \ y \ 1)^{\top} = H (u \ v \ 1)^{\top}$
 - \triangleright x, y are coordinates in the first plane
 - $lackbox{}{} u,v$ are coordinates in the second plane
- ▶ 9 elements but only 8 DoF, usually added constraint $h_{33} = 1$
- ► How to find H?
 - ► H, _ = cv2.findHomography(U, X)
 - ightharpoonup U, X are $N \times 2$ correspondence points

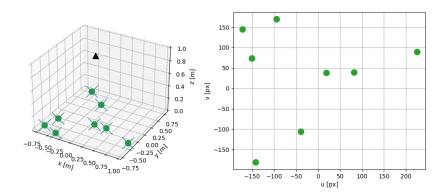
- lacktriangle Homography matrix H is 3×3 matrix that maps points from one plane to another
 - image plane to table desk
 - one image plane to another image plane (different view)
- - \triangleright x, y are coordinates in the first plane
 - $lackbox{}{} u,v$ are coordinates in the second plane
- ▶ 9 elements but only 8 DoF, usually added constraint $h_{33} = 1$
- ► How to find H?
 - ► H, _ = cv2.findHomography(U, X)
 - ightharpoonup 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

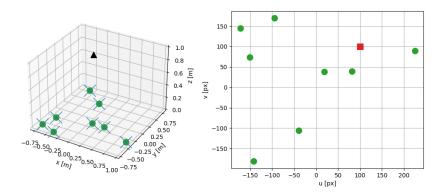


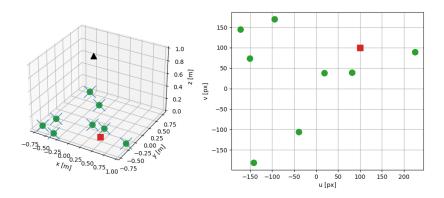








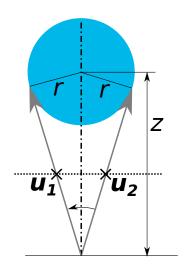




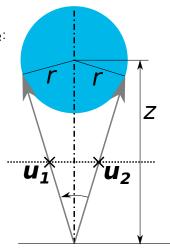
Non-planar pose estimation

- 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

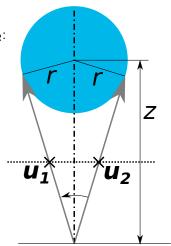
We know radius is fixed



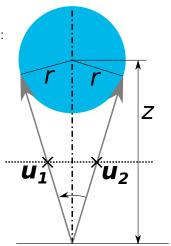
- We know radius is fixed
- From detected pixels u_1, u_2 , we can compute rays x_1, x_2 : $\frac{1}{\lambda_i}x_i = K^{-1}u_i$



- We know radius is fixed
- From detected pixels u_1, u_2 , we can compute rays x_1, x_2 : $\frac{1}{\lambda_i} x_i = K^{-1} u_i$
- Angle between vectors: $\cos \alpha = \frac{\frac{1}{\lambda_1 \lambda_2}}{\frac{1}{\lambda_1 \lambda_2}} \frac{\boldsymbol{x}_1 \cdot \boldsymbol{x}_2}{\|\boldsymbol{x}_1\| \|\boldsymbol{x}_2\|}$



- We know radius is fixed
- From detected pixels u_1,u_2 , we can compute rays x_1,x_2 : $\frac{1}{\lambda_i}x_i=K^{-1}u_i$
- ► Angle between vectors: $\cos \alpha = \frac{\frac{1}{\lambda_1 \lambda_2}}{\frac{1}{\lambda_1 \lambda_2}} \frac{x_1 \cdot x_2}{\|x_1\| \|x_2\|}$
- ▶ Depth: $z = \frac{r}{\sin(\alpha/2)}$



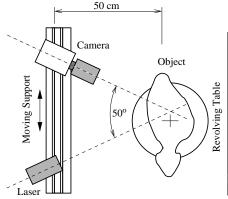
Using depth sensor

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

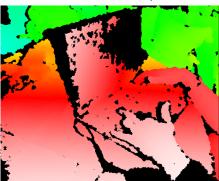




Issues with depth sensors

- ▶ 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)





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

Additional markers

► Can we compute the pose of patterns⁴?





⁴docs.opencv.org

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?







⁴docs.opencv.org

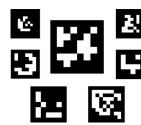


Additional markers

- ► 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



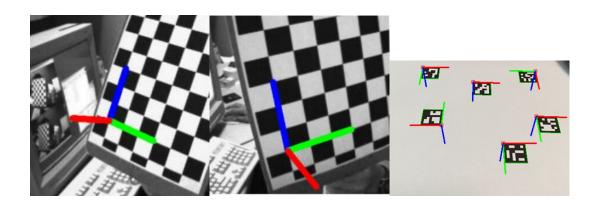




⁴docs.opencv.org



Markers pose example

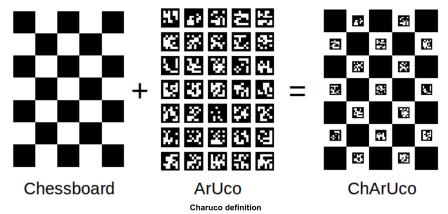


Markers example from real-world



ChArUco board for calibration

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



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

 \triangleright SE(3) poses of boards in camera frame

- 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

- 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?

- 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

- 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

- 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
- ► Mandatory excercise Safety
 - ► CIIRC, B670

What is next?

- ▶ Next week, we will start with a test
 - Questions from the first two lectures
 - ► SE2 and SE3 transformations and properties
 - Forward kinematics a DoF