

Introduction to Computer Vision

Lecture 1. Course Introduction

Sergey Zagoruyko

October 31, 2023

Table of Contents

Info about the course

Computer Vision problems & applications

Convolutions

Credits

- ▶ Mikhail Belyaev's slides
 - ▶ Founder of medical imaging solutions startup
 - ▶ AI for radiology <https://ira-labs.com/>
- ▶ <https://github.com/Lavton/SkoltechLaTeXtemplates>

Course info

Term 2, 3 credit course.

- ▶ **Tuesday, Thursday:** 9:00-12:00
 - ▶ Lectures: the first 90 minutes.
 - ▶ Practical exercises to work at home (provided after the lecture).
 - ▶ Online review of the exercises with QA after the next lecture: the remaining 60-90 minutes.
- ▶ 12 lectures + hands-ons.
- ▶ 1 invited lecture
- ▶ Final project.

Grading

- ▶ 3x20% - homeworks;
- ▶ 40% - final project.

Prerequisites

- ▶ **Basic knowledge of Python**
- ▶ Basic knowledge of linear algebra and statistics

Course outcomes

Knowledge

1. Statements of all major computer vision problems.
2. Mathematical details of the most important computer vision algorithms.

Skills

1. Select an appropriate method for solving particular computer vision problems.
2. Apply computer vision libraries.
3. Solve real-life computer vision problems.

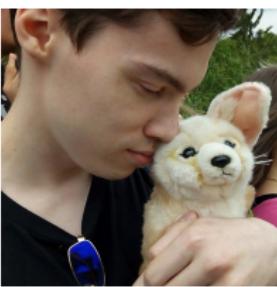
Knowledge

1. Implementing computer vision basic algorithms and complex pipelines.

Computer Vision Course



Sergey
Zagoruyko



Boris
Shirokikh



Maxim
Kurkin



Andrei
Filatov

Course structure

- ▶ Introduction to the field.
- ▶ Convolutions and Template Matching. Edges: Sobel filter, Canny
- ▶ 3D Reconstruction 1: Introduction and Keypoints
- ▶ 3D Reconstruction 2: Descriptors, Transforms, RANSAC
- ▶ 3D Reconstruction 3: Dense Reconstruction, SfM, Panorama Stitching
- ▶ Face Recognition - Metric Learning (invited lecture)
- ▶ Optical Flow
- ▶ Face Recognition
- ▶ Visual Tracking and Motion Prediction
- ▶ Computer Vision in Autonomous Vehicles
- ▶ Multimodal Computer Vision
- ▶ Neural Network Compression, Knowledge Distillation, Quantization

Course evaluation

Why do we need classical computer vision methods in 2023?

1. Many ideas in modern deep learning methods are based on classical computer vision methods.
2. Deep learning methods are data hungry. Using classical CV method you can create a prototype faster. Also, some unsupervised approaches can be used as to generate annotation candidates

Plagiarism Policy

Rules are simple:

1. If you use external code, please add the link
2. Homeworks are individual assignments. If one homework is submitted twice, both students get 1/2 of the score.
3. Final projects are team-based. However, do not submit a random project from the internet! It will hurt your overall score!

About me

- ▶ Started working with neural networks in 2010 at **MIPT**
- ▶ Graduated from **BMSTU** in 2013 with a degree in mechatronics.
- ▶ Defended a PhD in deep learning from **Ecole des Ponts ParisTech** in 2018.
- ▶ Completed internships at **FAIR**, **Intel**, and **VisionLabs**.
- ▶ Completed postdoctoral research at **Inria** and **FAIR**.
- ▶ Worked on autonomous vehicles at **Lyft** and **Toyota**. [my page](#)
- ▶ Currently serving as the head of fundamental research at **MTS AI** and assistant professor at **Skoltech**.



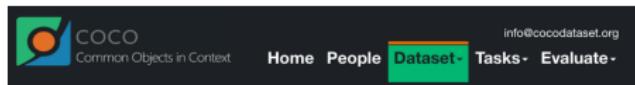
Computer vision in 2013



PhD in France in 2014



COCO dataset: internship at FAIR in 2015



COCO Explorer

COCO 2017 train/val browser (123,287 images, 886,284 instances). Crowd labels not shown.



elephant

2232 results

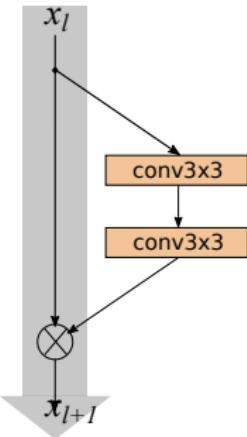


Try it out!

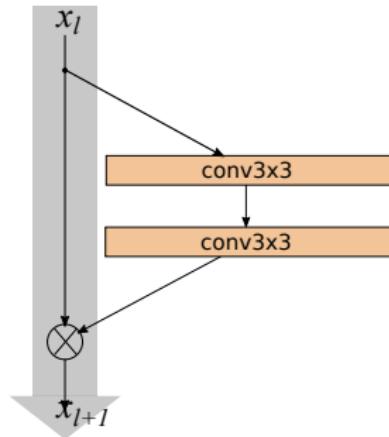


<https://cocodataset.org>

Wide Residual networks



ResNet



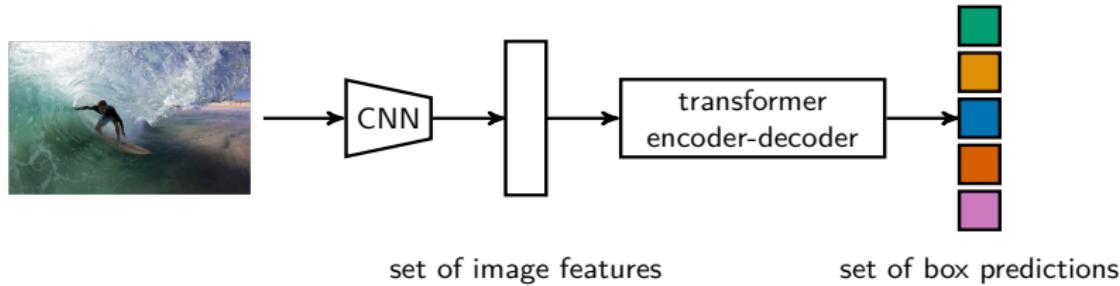
Wide ResNet

TITLE

CITED BY 7679 YEAR 2016

[Wide residual networks](#)
S Zagoruyko, N Komodakis
arXiv preprint arXiv:1605.07146

DETR



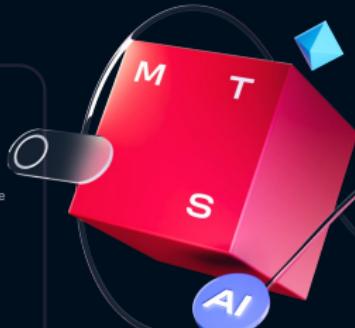
<input type="checkbox"/> TITLE	<input type="checkbox"/>	<input type="checkbox"/> ...	CITED BY	YEAR
<input type="checkbox"/> End-to-end object detection with transformers N Carion, F Massa, G Synnaeve, N Usunier, A Kirillov, S Zagoruyko European conference on computer vision, 213-229			8261	2020

About MTS AI

MTS AI ONE OF THE LEADING RUSSIAN COMPANIES IN AI DEVELOPMENT

20+ projects implemented

- Introduced chat bot to automate MTS customer service
- Using AI to predict profitability of investment in telecom infrastructure
- Taught KION to skip credits and locate commercials



Developing AI solutions and products

Based on CV, NLP and Edge Computing

Developing DeepTech and AI market

Investing in promising startups, helping corporations in AI-driven digitalization, organizing acceleration programs

R&D center

More than 200 experts and technological infrastructure: the most powerful supercomputer in telecommunications

MTS AI ASSETS



VisionLabs



INFOMOTIKO



Just AI



Kneron



Primo RPA

About the fundamental research team at MTS AI

- ▶ Natural Language Processing (NLP) Research:
 - ▶ Large Language Models (LLM)
 - ▶ Generative Pre-trained Transformer (GPT)
- ▶ Computer Vision (CV) Research:
 - ▶ Graph Reasoning
 - ▶ Multi-modal approaches combining image or video with text

Why do a PhD?

► Why do a PhD?

- ▶ Learn to formulate goals and plan research
- ▶ Learn to write articles
- ▶ Enjoy research freedom without product deadlines
- ▶ Interact with other students and graduate students
- ▶ Opportunity for a teaching career

► Why not do a PhD?

- ▶ Less computational resources compared to industry
- ▶ Lower salary during the period of the PhD
 - ▶ but possibly higher than the market after
- ▶ Intense competition in publications

Table of Contents

Info about the course

Computer Vision problems & applications

Convolutions

Computer Vision vs Image Processing

Computer Vision:

Involves interpreting and understanding images or videos.

- ▶ Input: Image or Video
- ▶ Output: Interpretation or understanding

Computer Vision vs Image Processing

Computer Vision:

Involves interpreting and understanding images or videos.

- ▶ Input: Image or Video
- ▶ Output: Interpretation or understanding

Image Processing:

Focuses on manipulating images or sets of images to enhance or extract information.

- ▶ Input: Image or Images
- ▶ Output: Modified Image or Images

Computer Vision vs Image Processing

Computer Vision:

Involves interpreting and understanding images or videos.

- ▶ Input: Image or Video
- ▶ Output: Interpretation or understanding

Image Processing:

Focuses on manipulating images or sets of images to enhance or extract information.

- ▶ Input: Image or Images
- ▶ Output: Modified Image or Images

Main Difference:

Computer Vision is about interpreting images, while Image Processing is about manipulating images to achieve a specific goal.

CV problems: Image Classification

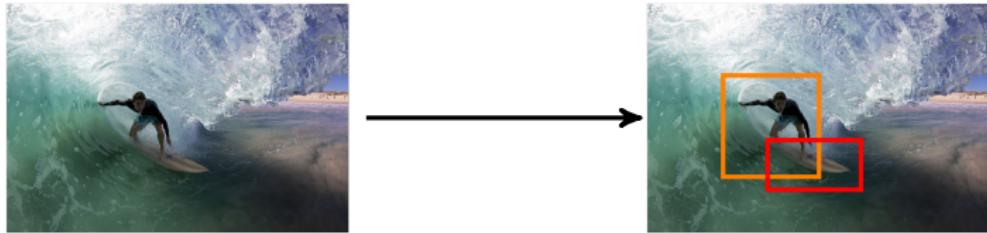


Image source: Krizhevsky et al. ImageNet Classification with Deep Convolutional Neural Networks

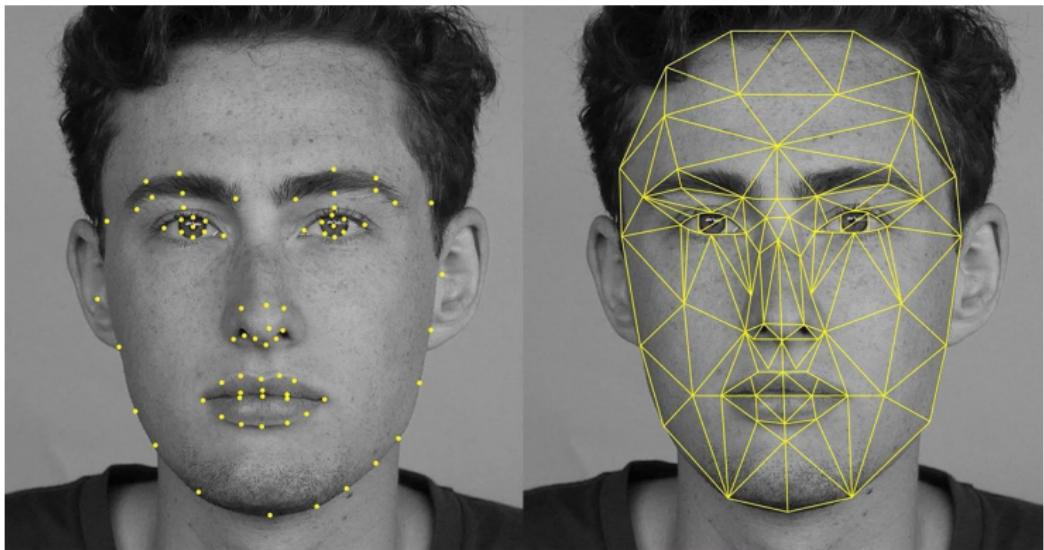
Object Detection Task

- Given an input image, predict **a set** of bounding boxes with corresponding categories

`set(("person", bbox1), ("surfboard", bbox2))`



CV problems: Keypoints Detection

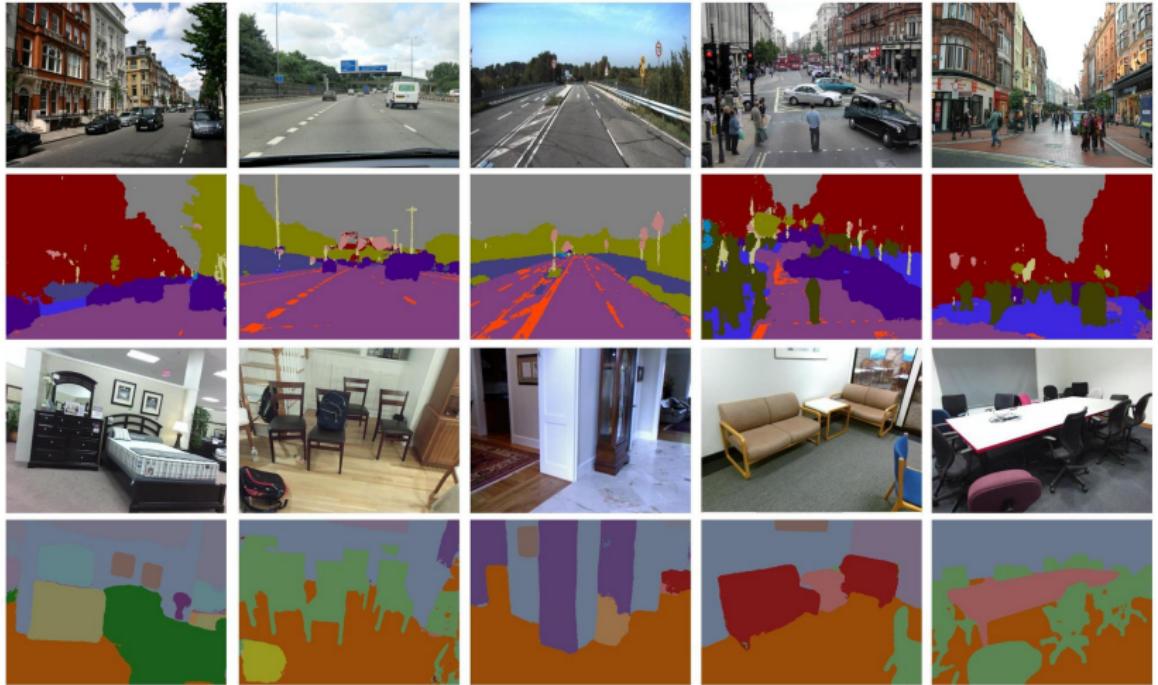


CV problems: Segmentation

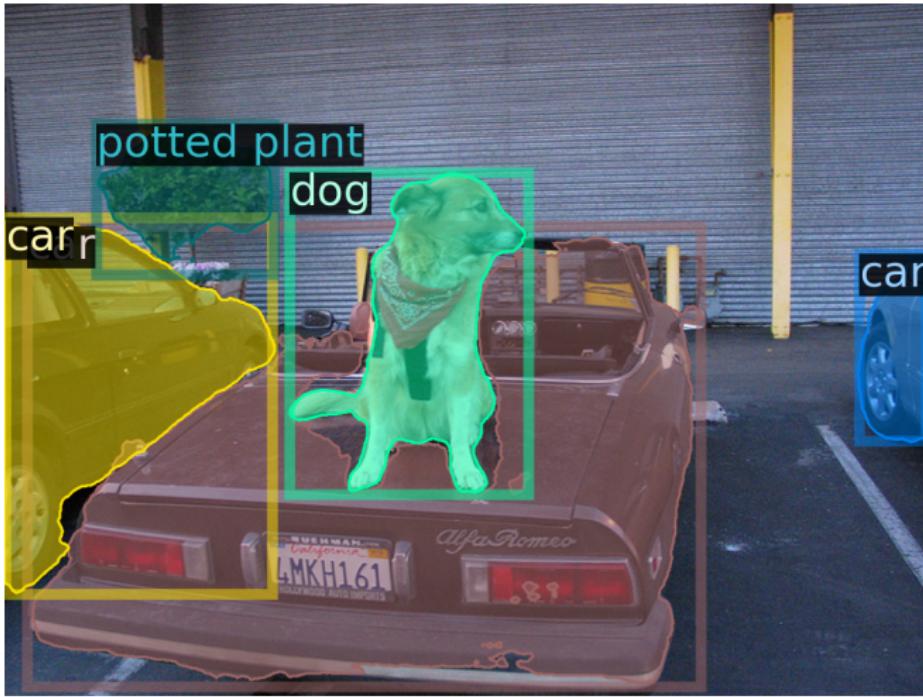


Image source: <https://www2.eecs.berkeley.edu/Research/Projects/CS/vision/grouping/resources.html>

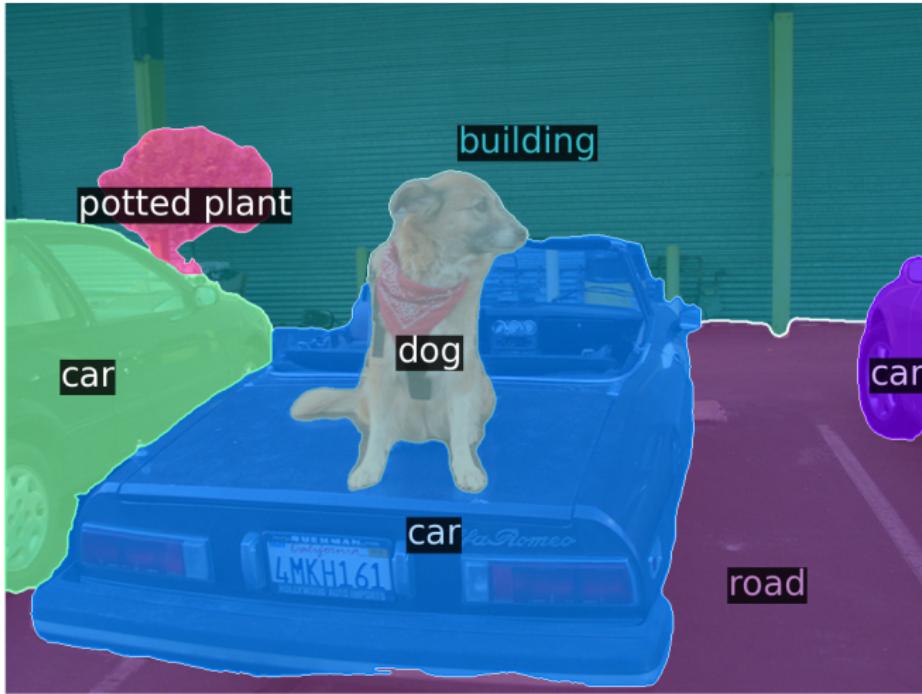
CV problems: Semantic Segmentation



CV problems: Instance Segmentation



CV problems: Panoptic Segmentation



CV problems: Image Captioning

A person riding a motorcycle on a dirt road.



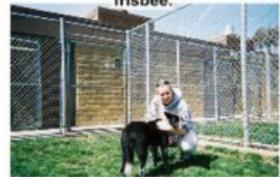
Two dogs play in the grass.



A skateboarder does a trick on a ramp.



A dog is jumping to catch a frisbee.



A group of young people playing a game of frisbee.



Two hockey players are fighting over the puck.



A little girl in a pink hat is blowing bubbles.



A refrigerator filled with lots of food and drinks.



A herd of elephants walking across a dry grass field.



A close up of a cat laying on a couch.



A red motorcycle parked on the side of the road.



A yellow school bus parked in a parking lot.



Image source: Vinyals et al. Show and Tell: A Neural Image Caption Generator

CV problems: Depth Estimation

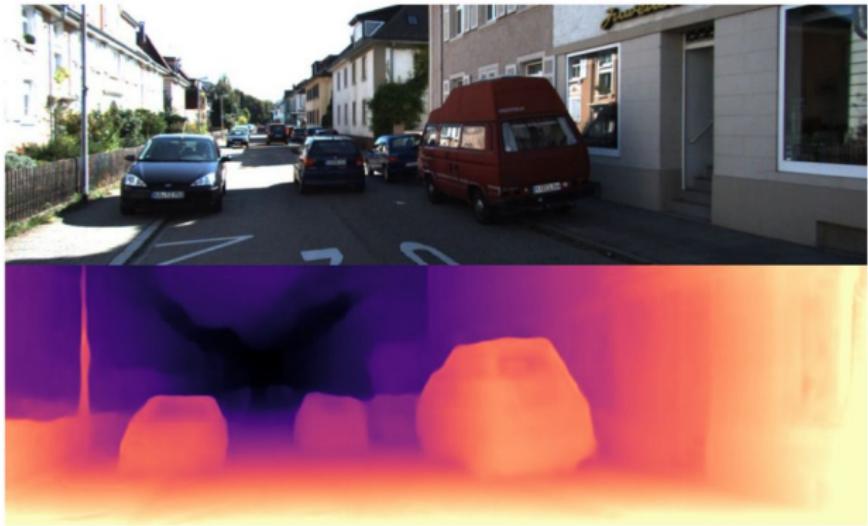
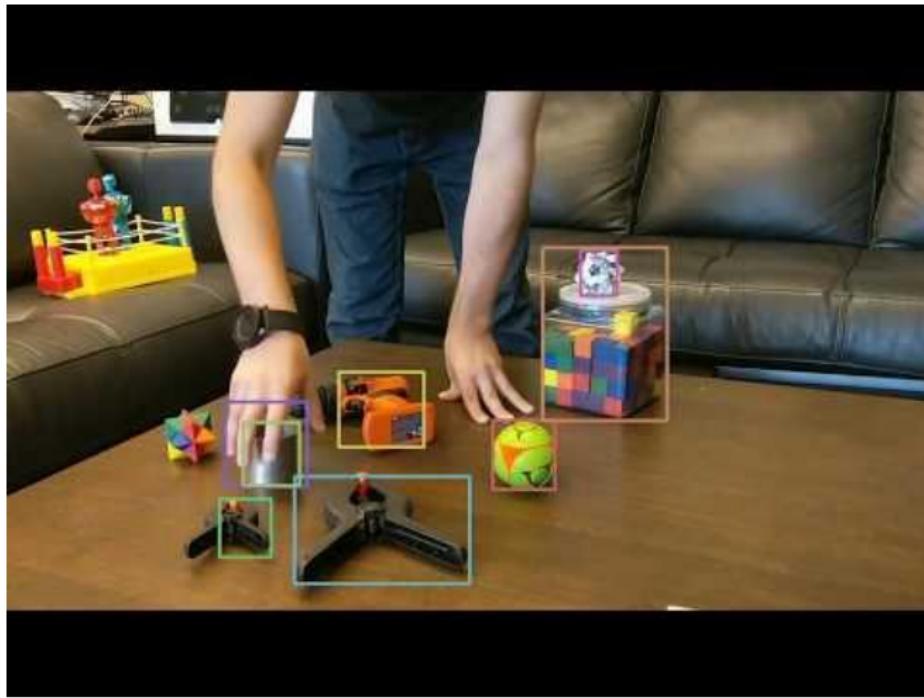


Image source: Vinyals et al. Show and Tell: A Neural Image Caption Generator

CV problems: Object Tracking



31 / 66 Video source: <https://arxiv.org/abs/1705.06368>

Image Processing problems: Super-resolution



Figure 2: From left to right: bicubic interpolation, deep residual network optimized for MSE, deep residual generative adversarial network optimized for a loss more sensitive to human perception, original HR image. Corresponding PSNR and SSIM are shown in brackets. [4× upscaling]

Image source: Ledig et al. Photo-Realistic Single Image Super-Resolution Using a Generative Adversarial Network

Image Processing problems: Inpainting

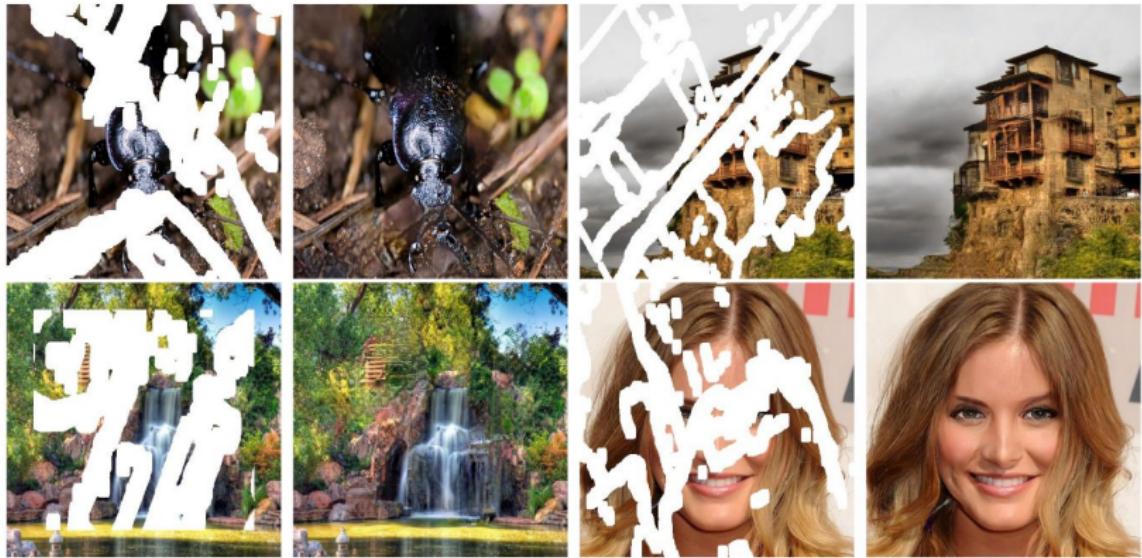


Image source: Liu et al. Image Inpainting for Irregular Holes Using Partial Convolutions

Image Processing problems: Image generation



Image sources: Gravity, Jurassic Park, Avatar

Image Processing problems: Style Transfer

Content image



Style image



Output image



+

+

+



Image Processing problems: Face generation

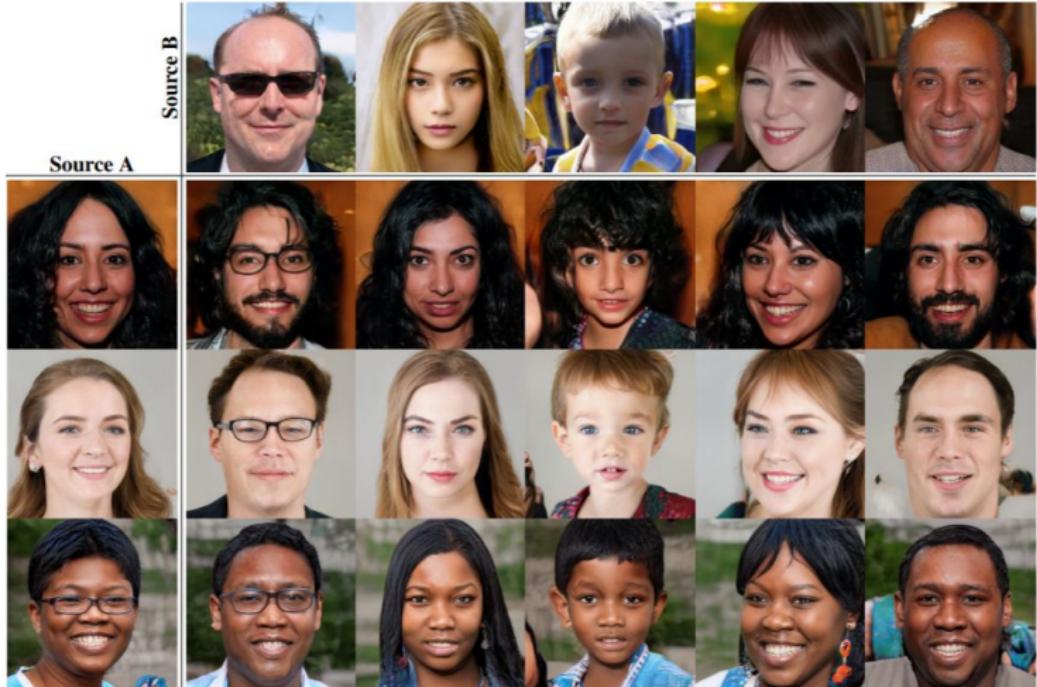
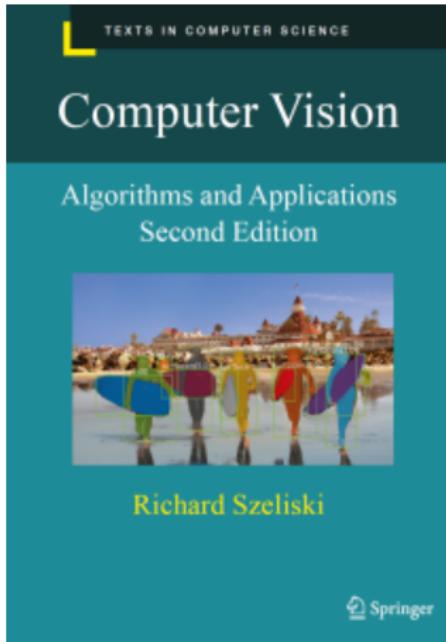


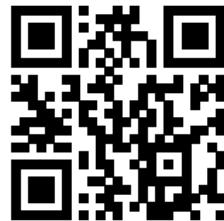
Image sources: <https://arxiv.org/abs/1812.04948>

Reading



The book is available online

Good reading for this lecture.



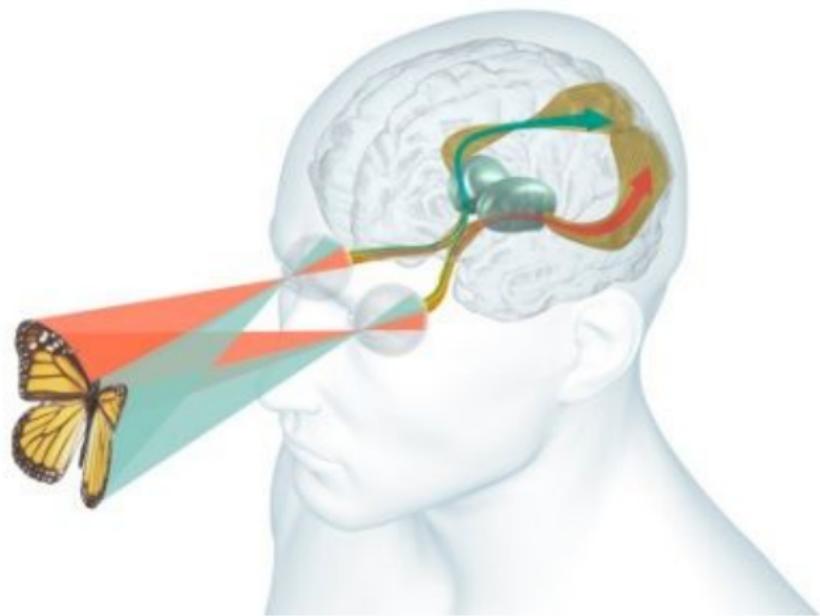
<https://szeliski.org/Book>

2nd edition

What is an image?



How do humans perceive images?



How do humans perceive images?



What are colors of this dress?

How do humans perceive images?

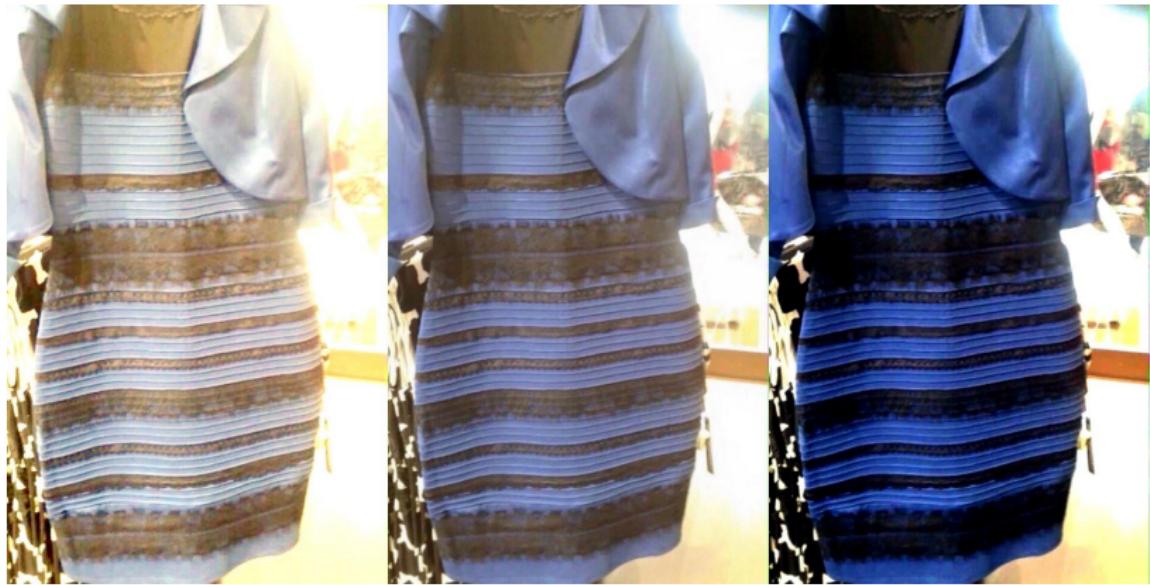


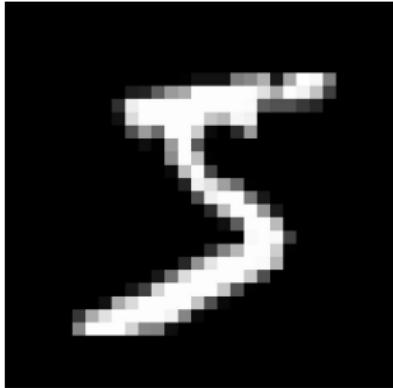
Image perception is subjective: different people will see this dress (left) either as blue and black or as white and brown.

How do humans perceive images?

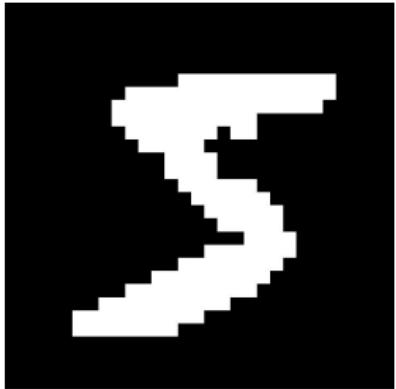


Image perception is subjective:
do you see a young woman or
an old one?

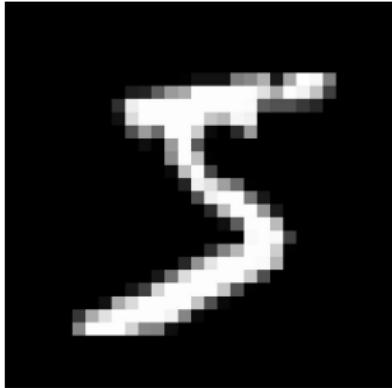
How do machines store images?



How do machines store images: Binary



How do machines store images: Grayscale



How do machines store images: RGB

Red



Green



How do machines store images: RGB

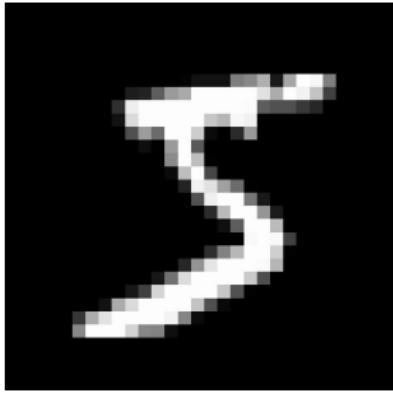
Red



Green



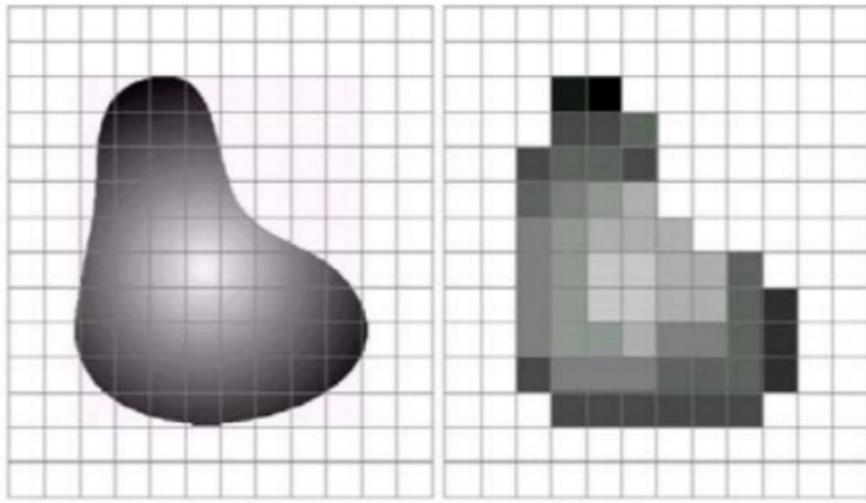
How do machines store images?



This a matrix. But from mathematical points of view we have another option.

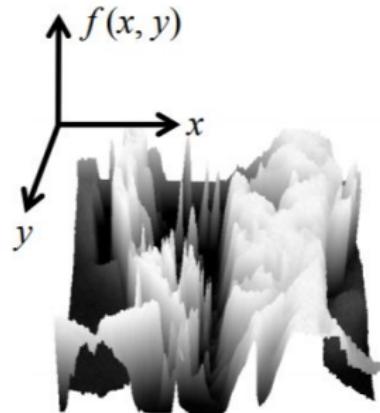
How do machines store images?

We can consider image as a discrete representation of a 2D function.



How do machines store images?

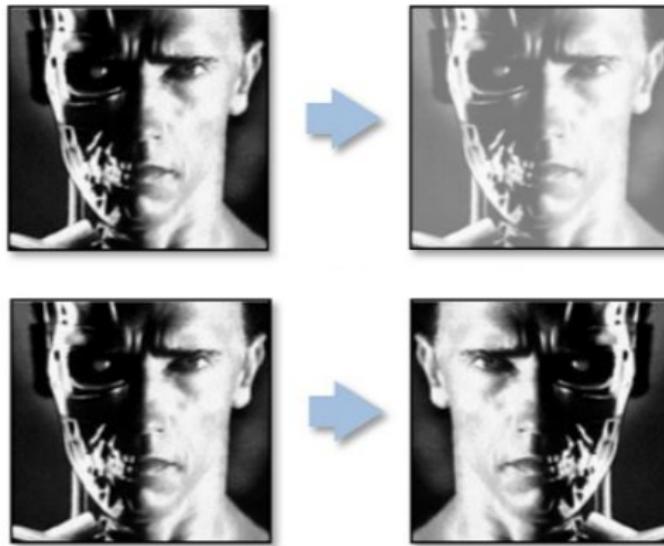
We can consider image as a discrete representation of a 2D function.



Source: N. Snavely

How do machines store images?

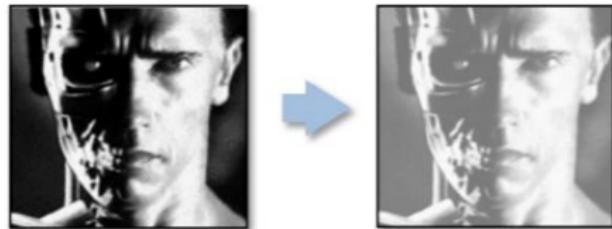
So some standard operations can be applied



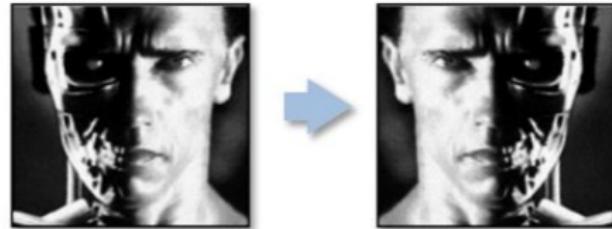
Source: N. Snavely

How do machines store images?

So some standard operations can be applied



$$g(x,y) = f(x,y) + 20$$



$$g(x,y) = f(-x,y)$$

Table of Contents

Info about the course

Computer Vision problems & applications

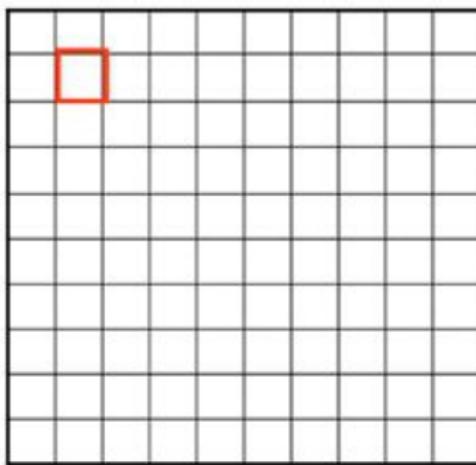
Convolutions

Convolutions: moving average example

$F[x, y]$

0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	90	90	90	90	90	0
0	0	0	90	90	90	90	90	90	0
0	0	0	90	90	90	90	90	90	0
0	0	0	90	0	90	90	90	90	0
0	0	0	90	90	90	90	90	90	0
0	0	0	0	0	0	0	0	0	0
0	0	0	90	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

$G[x, y]$



Source: S. Seitz

Convolutions: moving average example

$F[x, y]$

0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0
0	0	0	90	90	90	90	90	0	0	0
0	0	0	90	90	90	90	90	0	0	0
0	0	0	90	90	90	90	90	0	0	0
0	0	0	90	0	90	90	90	0	0	0
0	0	0	90	90	90	90	90	0	0	0
0	0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0

$G[x, y]$

	0	10	20	30	30	30	20	10		
	0	20	40	60	60	60	40	20		
	0	30	60	90	90	90	60	30		
	0	30	50	80	80	90	60	30		
	0	30	50	80	80	90	60	30		
	0	20	30	50	50	60	40	20		
10	20	30	30	30	30	30	20	10		
10	10	10	0	0	0	0	0	0		

Source: S. Seitz

Convolutions: some math

Let us start with general 2 dimensional functions

$$(f * g)(x, y) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x', y') g(x - x', y - y') dx' dy' \quad (1)$$

Convolutions: some math

Let us start with general 2 dimensional functions

$$(f * g)(x, y) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x', y') g(x - x', y - y') dx' dy' \quad (1)$$

But images are discrete representations of 2D functions

$$(f * g)(x, y) = \sum_{x'=-\infty}^{\infty} \sum_{y'=-\infty}^{\infty} f(x', y') g(x - x', y - y') \quad (2)$$

Convolutions: some math

Let us start with general 2 dimensional functions

$$(f * g)(x, y) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x', y') g(x - x', y - y') dx' dy' \quad (1)$$

But images are discrete representations of 2D functions

$$(f * g)(x, y) = \sum_{x'=-\infty}^{\infty} \sum_{y'=-\infty}^{\infty} f(x', y') g(x - x', y - y') \quad (2)$$

And, finally, images have finite support

$$(f * g)(x, y) = \sum_{x'=0}^{X} \sum_{y'=0}^{Y} f(x', y') g(x - x', y - y') \quad (3)$$

Convolutions: some math

Question: what is g for moving average?

$$(f * g)(x, y) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x', y') g(x - x', y - y') dx' dy' \quad (4)$$

Convolutions: some math

Question: what is g for moving average?

$$(f * g)(x, y) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x', y') g(x - x', y - y') dx' dy' \quad (4)$$

$$\begin{bmatrix} \frac{1}{9} & \frac{1}{9} & \frac{1}{9} \\ \frac{1}{9} & \frac{1}{9} & \frac{1}{9} \\ \frac{1}{9} & \frac{1}{9} & \frac{1}{9} \end{bmatrix}$$

$$\frac{1}{9} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$$

This operation & the corresponding matrix are also called filters, kernels, convolutional matrices.

Convolutions: moving average

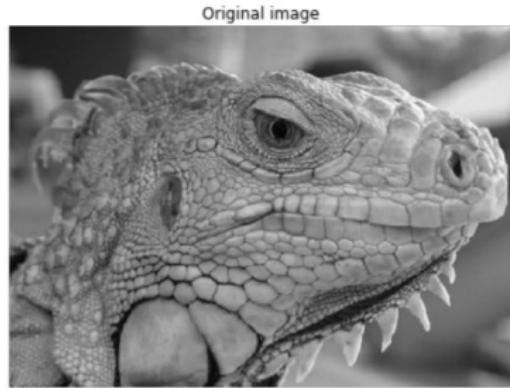
How does it modify images?

$$\frac{1}{9} \begin{array}{|c|c|c|} \hline 1 & 1 & 1 \\ \hline 1 & 1 & 1 \\ \hline 1 & 1 & 1 \\ \hline \end{array}$$

Convolutions: moving average

How does it modify images?

$$\frac{1}{9} \begin{array}{|c|c|c|} \hline 1 & 1 & 1 \\ \hline 1 & 1 & 1 \\ \hline 1 & 1 & 1 \\ \hline \end{array}$$



Convolutions in 2D

1 <small>$\times 1$</small>	1 <small>$\times 0$</small>	1 <small>$\times 1$</small>	0	0
0 <small>$\times 0$</small>	1 <small>$\times 1$</small>	1 <small>$\times 0$</small>	1	0
0 <small>$\times 1$</small>	0 <small>$\times 0$</small>	1 <small>$\times 1$</small>	1	1
0	0	1	1	0
0	1	1	0	0

Image

4		

Convolved
Feature

Convolutions in 2D: Practical exercise 1

Let's implement a Python function that convolves an image with a kernel.

$$(f * g)(x, y) = \sum_{x'=0}^X \sum_{y'=0}^Y f(x', y')g(x - x', y - y') \quad (5)$$

Convolutions: examples of filters

How does it modify images?

•0	•0	•0
•0	•1	•0
•0	•0	•0

Convolutions: examples of filters

How does it modify images?

$$\begin{matrix} \bullet 0 & \bullet 0 & \bullet 0 \\ \bullet 0 & \bullet 1 & \bullet 0 \\ \bullet 0 & \bullet 0 & \bullet 0 \end{matrix}$$



*

$$\begin{matrix} \bullet 0 & \bullet 0 & \bullet 0 \\ \bullet 0 & \bullet 1 & \bullet 0 \\ \bullet 0 & \bullet 0 & \bullet 0 \end{matrix}$$

=



Convolutions: examples of filters

How does it modify images?

•0	•0	•0
•0	•0	•1
•0	•0	•0

Convolutions: examples of filters

How does it modify images?

•0	•0	•0
•0	•0	•1
•0	•0	•0



*

•0	•0	•0
•0	•0	•1
•0	•0	•0

=



Original

Shifted right
By 1 pixel

Convolutions: examples of filters

How does it
modify
images?

$$\begin{array}{|c|c|c|} \hline \bullet 0 & \bullet 0 & \bullet 0 \\ \hline \bullet 0 & \bullet 2 & \bullet 0 \\ \hline \bullet 0 & \bullet 0 & \bullet 0 \\ \hline \end{array} - \frac{1}{9} \begin{array}{|c|c|c|} \hline \bullet 1 & \bullet 1 & \bullet 1 \\ \hline \bullet 1 & \bullet 1 & \bullet 1 \\ \hline \bullet 1 & \bullet 1 & \bullet 1 \\ \hline \end{array} = ?$$

Convolutions: examples of filters



Original

$$\begin{bmatrix} \bullet 0 & \bullet 0 & \bullet 0 \\ \bullet 0 & \bullet 2 & \bullet 0 \\ \bullet 0 & \bullet 0 & \bullet 0 \end{bmatrix}$$

-

$$\frac{1}{9} \begin{bmatrix} \bullet 1 & \bullet 1 & \bullet 1 \\ \bullet 1 & \bullet 1 & \bullet 1 \\ \bullet 1 & \bullet 1 & \bullet 1 \end{bmatrix}$$

= ?

(Note that filter sums to 1)

“details of the image”

$$\begin{bmatrix} \bullet 0 & \bullet 0 & \bullet 0 \\ \bullet 0 & \bullet 1 & \bullet 0 \\ \bullet 0 & \bullet 0 & \bullet 0 \end{bmatrix}$$

+

$$\begin{bmatrix} \bullet 0 & \bullet 0 & \bullet 0 \\ \bullet 0 & \bullet 1 & \bullet 0 \\ \bullet 0 & \bullet 0 & \bullet 0 \end{bmatrix}$$

-

$$\frac{1}{9} \begin{bmatrix} \bullet 1 & \bullet 1 & \bullet 1 \\ \bullet 1 & \bullet 1 & \bullet 1 \\ \bullet 1 & \bullet 1 & \bullet 1 \end{bmatrix}$$

Source: J. Niebles

Convolutions: examples of filters



- Let's add it back:



Convolutions: examples of filters



Original

$$\begin{bmatrix} \bullet 0 & \bullet 0 & \bullet 0 \\ \bullet 0 & \bullet 2 & \bullet 0 \\ \bullet 0 & \bullet 0 & \bullet 0 \end{bmatrix}$$

-

$$\frac{1}{9}$$

$$\begin{bmatrix} \bullet 1 & \bullet 1 & \bullet 1 \\ \bullet 1 & \bullet 1 & \bullet 1 \\ \bullet 1 & \bullet 1 & \bullet 1 \end{bmatrix}$$

= ?

(Note that filter sums to 1)



Original

$$\begin{bmatrix} \bullet 0 & \bullet 0 & \bullet 0 \\ \bullet 0 & \bullet 2 & \bullet 0 \\ \bullet 0 & \bullet 0 & \bullet 0 \end{bmatrix}$$

-

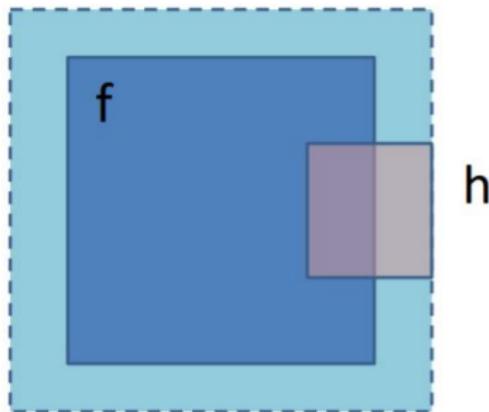
$$\frac{1}{9}$$

$$\begin{bmatrix} \bullet 1 & \bullet 1 & \bullet 1 \\ \bullet 1 & \bullet 1 & \bullet 1 \\ \bullet 1 & \bullet 1 & \bullet 1 \end{bmatrix}$$

=



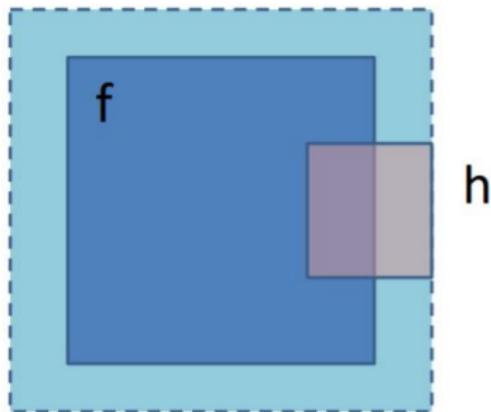
Convolutions: boundary conditions



How to process image boundaries?

Source: J. Niebles

Convolutions: boundary conditions



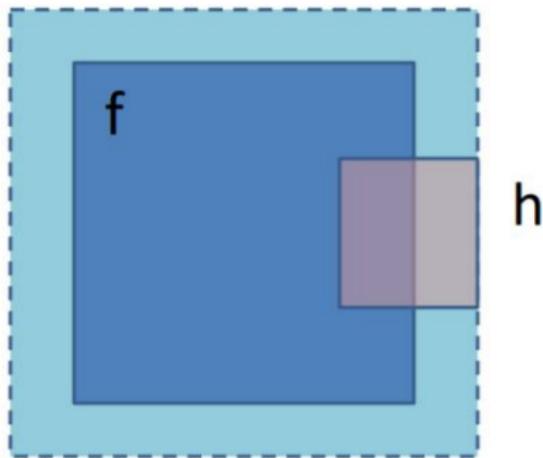
How to process image boundaries?

- ▶ **Reduce size** of the processed image
- ▶ **Add zeros** or constants
- ▶ Mirror the image

Source: J. Niebles

Convolutions: boundary conditions - Task 2

Let's extend our convolution function with simple zero padding.



Source: J. Niebles

Welcome!



Course telegram channel