# Object Recognition/Detection
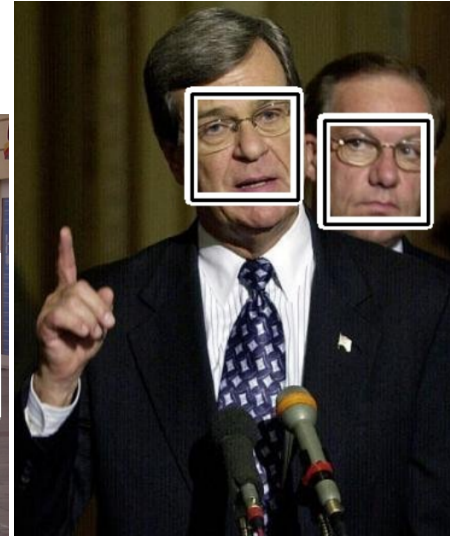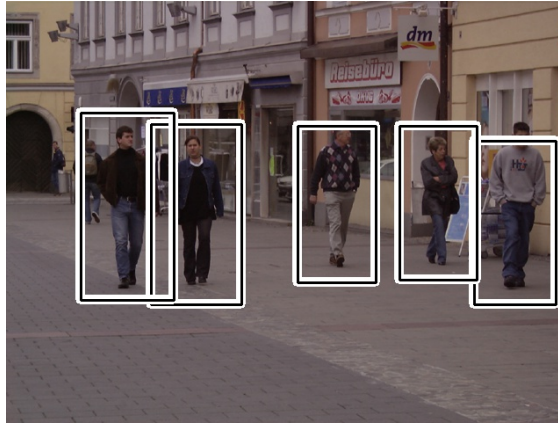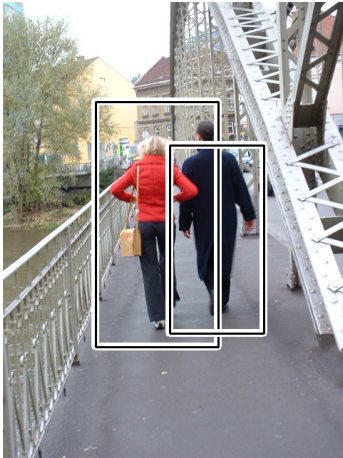
# Radovan Fusek

2nd International summer school on "Deep Learning and Visual Data Analysis"

2018

# What is Object Detection/Recognition?

- Output?
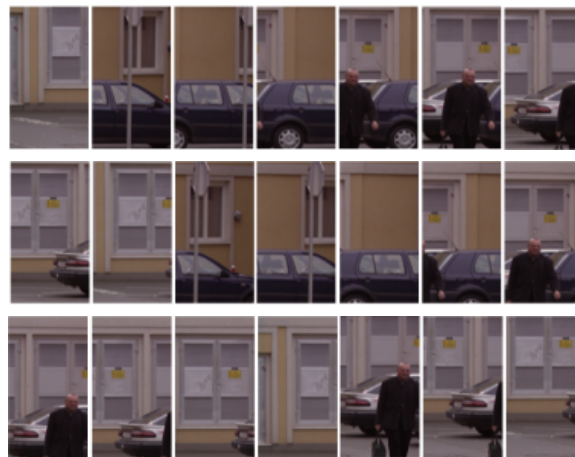    - position of the objects
    - scale of the objects
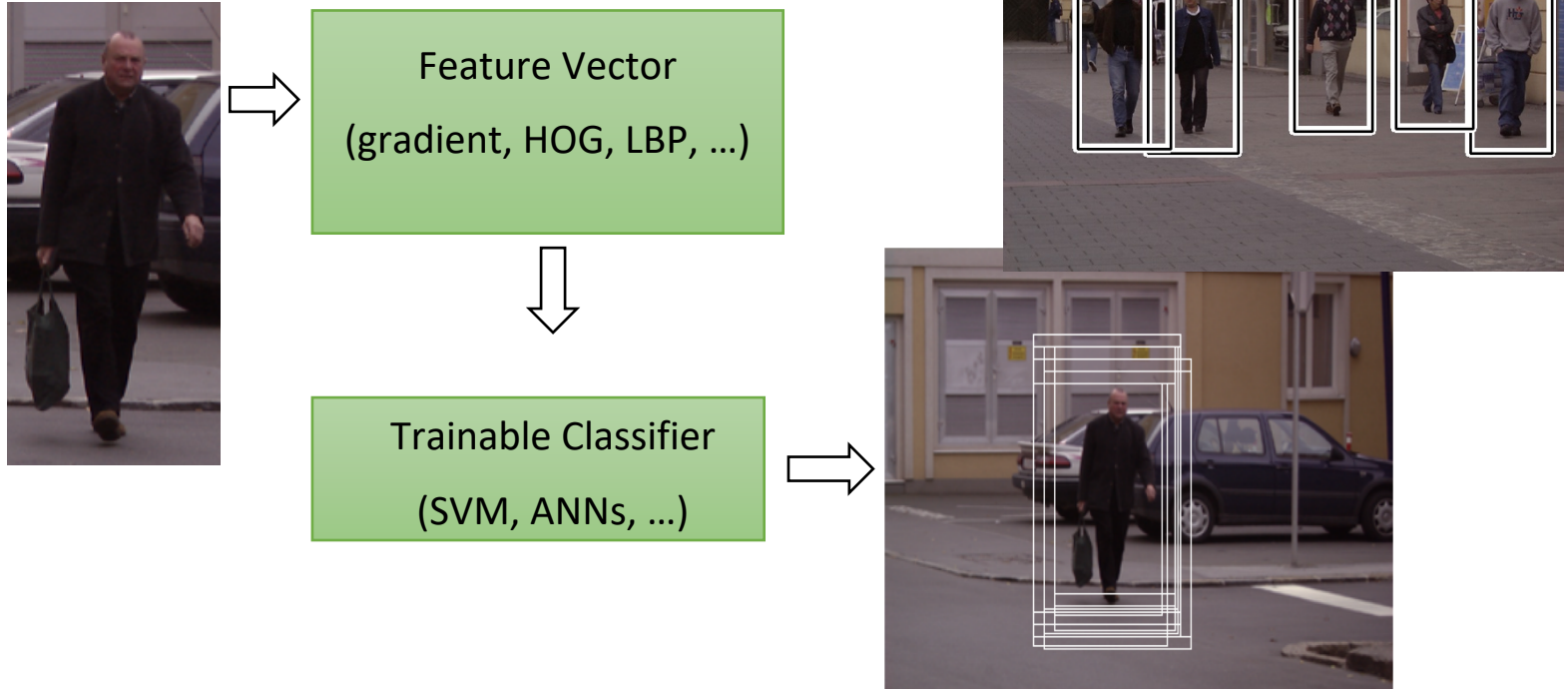    - name of the objects

# Object Detection/Recognition

- Haar

- HOG    } Traditional Approaches

- LBP

- SIFT, SURF    } KeyPoints

- CNNs    } Deep Learning Approach

- Practical examples using OpenCV + Dlib (https://opencv.org/, http://dlib.net/)

# Sliding Window - Main Idea

Constantine Papageorgiou and Tomaso Poggio: A Trainable System for Object Detection.
*Int. J. Comput. Vision* 38, pp. 15-33. (2000)

# Related Works



Feature Vector

(gradient, HOG, LBP, …)

Trainable Classifier

(SVM, ANNs, …)

Constantine Papageorgiou and Tomaso Poggio: A Trainable System for Object Detection.
*Int. J. Comput. Vision* 38, pp. 15-33. (2000)
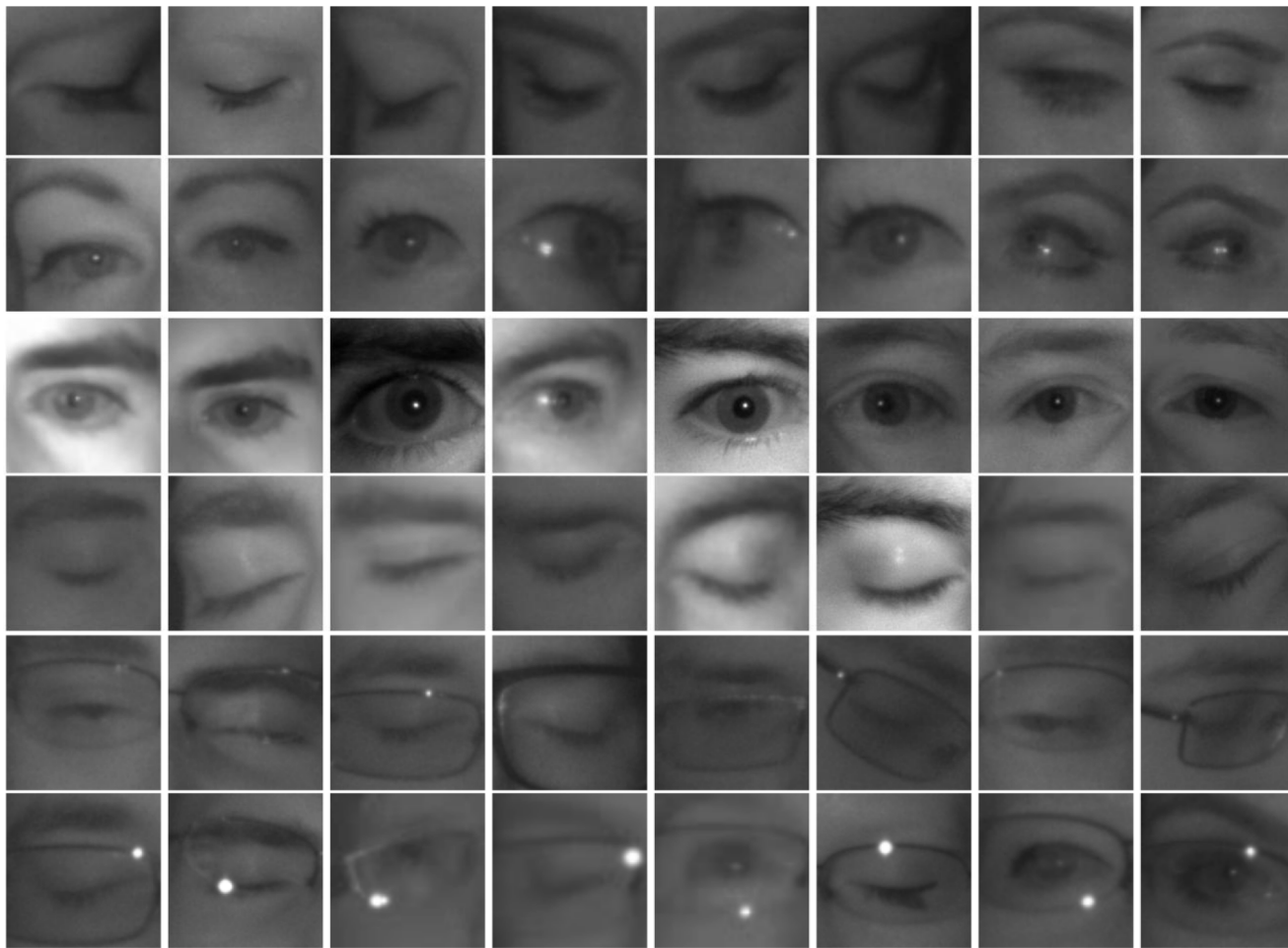
# Generating Training Set

- negative set - without the object of interest
- positive set
  - rotation
  - noise
  - Illumination
  - scale

# Generating Training Set

# Object Detection/Recognition

- <span style="color:red">Haar</span>

- HOG                    <span style="color:red">Traditional Approaches</span>

- LBP

- SIFT, SURF            KeyPoints

- CNNs                   Deep Learning Approach

- Practical examples using OpenCV + Dlib (https://opencv.org/, http://dlib.net/)

# Related Works

# Features

- faces have similar properties
    - eye regions are darker than the upper-cheeks
    - the nose bridge region is brighter than the eyes



https://docs.opencv.org/3.4.1/d7/d8b/tutorial_py_face_detection.html

# Features

- Rectangular features

$$F_{Haar} = E(R_{white}) - E(R_{black})$$

$R_{white}$   $R_{black}$

# Features

**Different sets**



Basic Haar set

Haar-like features

Extended Haar features

Borders

Lines

Center-surround

# Feature Selection

# Feature Selection

- AdaBoost (Adaptive Boost) is an iterative learning algorithm to construct a "strong" classifier as a linear combination of weighted simple "weak" classifiers

- weak classifier - each single rectangle feature (features as weak classifiers)

- during each iteration, each example/image receives a weight determining its importance

# Feature Selection

p AdaBoost starts with a uniform distribution of "weights" over training examples.

p Select the classifier with the lowest weighted error (i.e. a "weak" classifier)

p Increase the weights on the training examples that were misclassified.

p (Repeat)

p At the end, carefully make a linear combination of the weak classifiers obtained at all iterations.

# Cascade of Classifier

The idea of cascade classifier is reject the non-face region as soon as possible

# Cascade of Classifier

The idea of cascade classifier is reject the non-face region as soon as possible

# Cascade of Classifier

The idea of cascade classifier is reject the non-face region as soon as possible

# Cascade of Classifier

The idea of cascade classifier is reject the non-face region as soon as possible

# Cascade of Classifier

The idea of cascade classifier is reject the non-face region as soon as possible

# Cascade of Classifier

The idea of cascade classifier is reject the non-face region as soon as possible

# Cascade of Classifier

The idea of cascade classifier is reject the non-face region as soon as possible
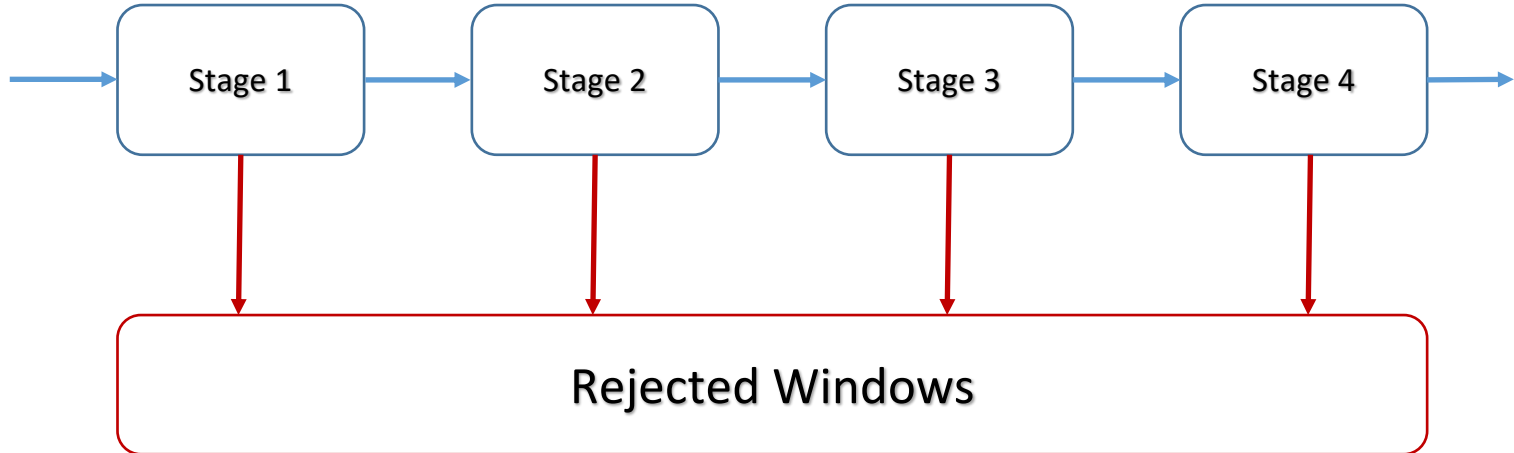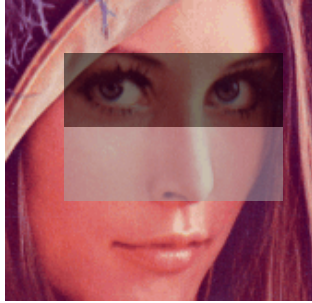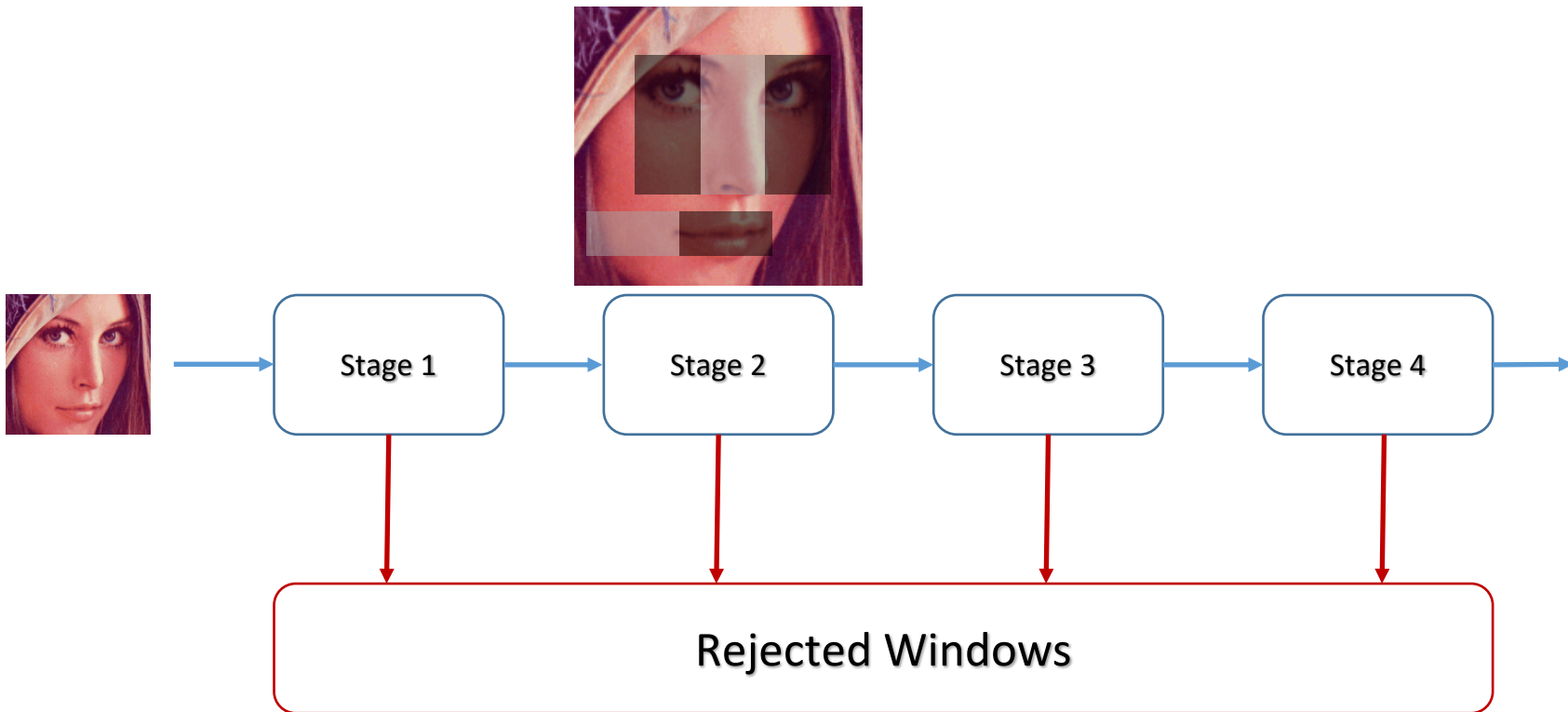
# Cascade of Classifier

The idea of cascade classifier is reject the non-face region as soon as possible

# Haar Features



https://vimeo.com/12774628

# Parking Lot Occupation

- Fabián, T.: **A Vision-based Algorithm for Parking Lot Utilization Evaluation Using Conditional Random Fields**. In 9th International Symposium on Visual Computing ISVC 2013, pp. 1-12 (2013)

- Fusek, R., Mozdřeň, K., Šurkala, M., Sojka, E.: **AdaBoost for Parking Lot Occupation Detection**. Advances in Intelligent Systems and Computing, vol. 226, pp. 681-690 (2013)

**http://mrl.cs.vsb.cz/**

# Haar Features

The modified version of Haar-like features that more properly reflect the shape of the pedestrians than the classical Haar-like features.

Hoang, V.D., Vavilin, A., Jo, K.H.: Pedestrian detection approach based on modified haar-like features and adaboost. In: Control, Automation and Systems (ICCAS), 2012 12th International Conference on. pp. 614-618 (Oct 2012)

# Object Detection/Recognition

- Haar

- HOG

- LBP

Traditional Approaches

- SIFT, SURF

KeyPoints

- CNNs

Deep Learning Approach

- Practical examples using OpenCV + Dlib (https://opencv.org/, http://dlib.net/)

# Related Works

Papageorgiou
(2000)

**A Trainable System for Object Detection**

CONSTANTINE PAPAGEORGIOU AND TOMASO POGGIO
*Center for Biological and Computational Learning, Artificial Intelligence Laboratory, MIT,
Cambridge, MA, USA*
cpapa@ai.mit.edu
tp@ai.mit.edu



Viola, Jones
(2001,2004)

**Robust Real-Time Face Detection**

PAUL VIOLA
*Microsoft Research, One Microsoft Way, Redmond, WA 98052, USA*
viola@microsoft.com

MICHAEL J. JONES
*Mitsubishi Electric Research Laboratory, 201 Broadway, Cambridge, MA 02139, USA*
mjones@merl.com



Dalal, Triggs
(2005)
cit. 10947

**Histograms of Oriented Gradients for Human Detection**

**Navneet Dalal and Bill Triggs**
INRIA Rhône-Alps, 655 avenue de l'Europe, Montbonnot 38334, France
{Navneet.Dalal,Bill.Triggs}@inrialpes.fr,   http://lear.inrialpes.fr

# Histograms of Oriented Gradients (HOG)

**Basic Steps:**

- In HOG, a sliding window is used for detection.

- The window is divided into small connected cells.

- The histograms of gradient orientations are calculated in each cell.

- Support Vector Machine (SVM) classifier.



Input image

Detection window

**Normalise gamma**

**Compute gradients**

**Weighted vote in spatial & orientation cells**

**Contrast normalise over overlapping spatial cells**

**Collect HOGs over detection window**

**Linear SVM**

Cell

Block

Overlap of Blocks

Feature vector $f =$
[ ..., ...,    ...]

# Histograms of Oriented Gradients (HOG)

**Blocks, Cells:**

# Histograms of Oriented Gradients (HOG)



**Blocks, Cells:**

- 8 x 8 cell

- 16 x 16 block – overlap

- normalization within the blocks

**Final Vector: Collect HOG blocks into vector**

# Histograms of Oriented Gradients (HOG)

# Practical Example – Detection + Recognition

Consider the following problem: Find and recognize two following lego kits

# OpenCV - http://opencv.org/

## OpenCV 3.3.0
### Open Source Computer Vision

| Main Page | Related Pages | Modules | Namespaces ▾ | Classes ▾ | Files ▾ | Examples |

## Introduction

OpenCV (Open Source Computer Vision Library: http://opencv.org) is an open-source BSD-licensed library that includes several hundreds of computer vision algorithms. The document describes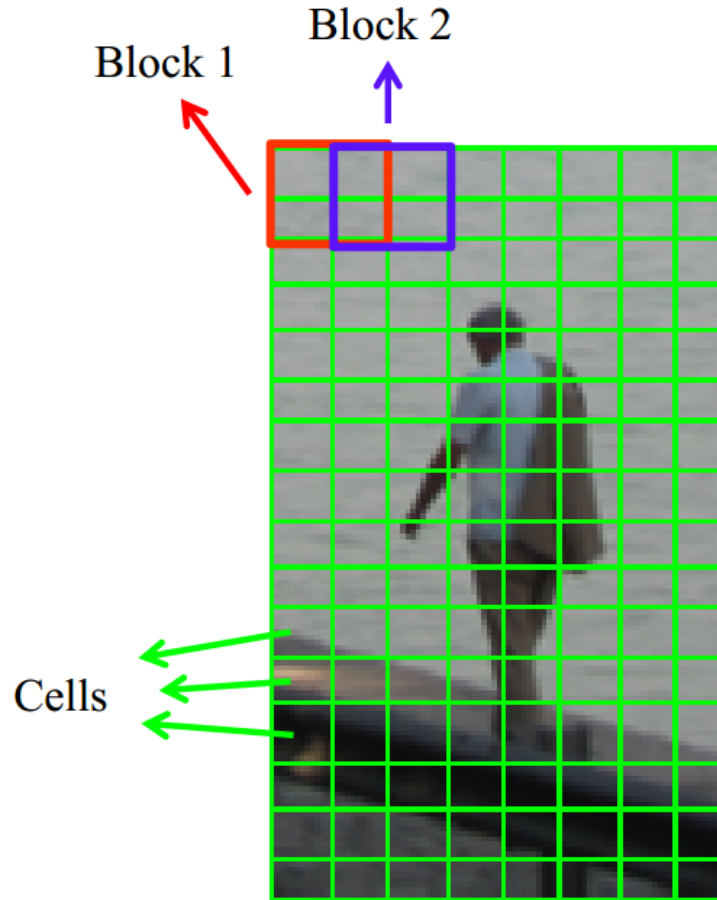 the so-called OpenCV 2.x API, which is essentially a C++ API, as opposite to the C-based OpenCV 1.x API. The latter is described in opencv1x.pdf.
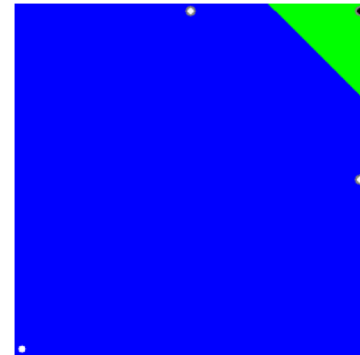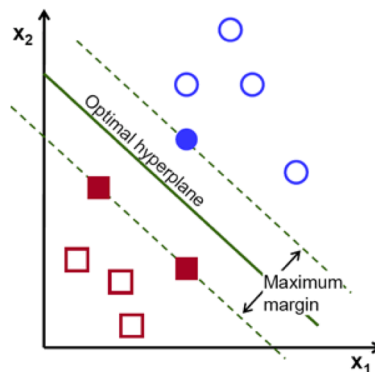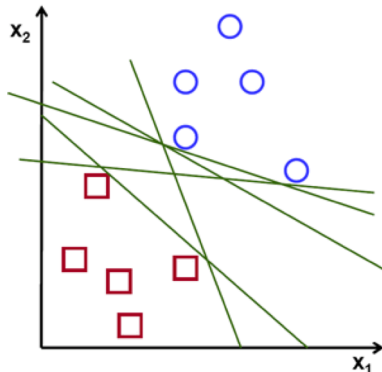
OpenCV has a modular structure, which means that the package includes several shared or static libraries. The following modules are available:

- **Core functionality** - a compact module defining basic data structures, including the dense multi-dimensional array Mat and basic functions used by all other modules.
- **Image processing** - an image processing module that includes linear and non-linear image filtering, geometrical image transformations (resize, affine and perspective warping, generic table-based remapping), color space conversion, histograms, and so on.
- **video** - a video analysis module that includes motion estimation, background subtraction, and object tracking algorithms.
- **calib3d** - basic multiple-view geometry algorithms, single and stereo camera calibration, object pose estimation, stereo correspondence algorithms, and elements of 3D reconstruction.
- **features2d** - salient feature detectors, descriptors, and descriptor matchers.
- **objdetect** - detection of objects and instances of the predefined classes (for example, faces, eyes, mugs, people, cars, and so on).
- **highgui** - an easy-to-use interface to simple UI capabilities.
- **Video I/O** - an easy-to-use interface to video capturing and video codecs.
- **gpu** - GPU-accelerated algorithms from different OpenCV modules.
- ... some other helper modules, such as FLANN and Google test wrappers, Python bindings, and others.
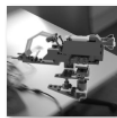
The further chapters of the document describe functionality of each module. But first, make sure to get familiar with the common API concepts used thoroughly in the library.

# Detection step - HOG+SVM (OpenCV)

```cpp
1    // Set up training data
2    int labels[4] = {1, -1, -1, -1};
3    Mat labelsMat(4, 1, CV_32SC1, labels);
4
5    float trainingData[4][2] = { {501, 10}, {255, 10}, {501, 255}, {10, 501} };
6    Mat trainingDataMat(4, 2, CV_32FC1, trainingData);
7
8    // Set up SVM's parameters
9    SVM::Params params;
10   params.svmType    = SVM::C_SVC;
11   params.kernelType = SVM::LINEAR;
12   params.termCrit   = TermCriteria(TermCriteria::MAX_ITER, 100, 1e-6);
13
14   // Train the SVM
15   Ptr<SVM> svm = StatModel::train<SVM>(trainingDataMat, ROW_SAMPLE, labelsMat, params);
```
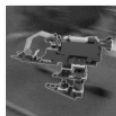


https://docs.opencv.org/3.1.0/d1/d73/tutorial_introduction_to_svm.html

# Alien

# Avenger

# Detection step - HOG+SVM (OpenCV)

# Detection step - HOG+SVM (OpenCV)
# Sliding Window (detectMultiScale)

```cpp
1    int blockSize = 16;
2    int cellSize = 8;
3    int strideSize = 8;
4    int winSize = 64;
5
6    //HOGDescriptor hog;
7    HOGDescriptor my_hog(
8              cv::Size(winSize,winSize), //winSize
9              cv::Size(blockSize,blockSize), //blocksize
10             cv::Size(strideSize,strideSize), //blockStride,
11             cv::Size(cellSize,cellSize), //cellSize,
12             9, //nbins,
13             );
14
15    //SVM
16    Ptr<SVM> svm = StatModel::load<SVM>( classifierName );
17    std::vector< float > hog_detector;
18    //get the support vectors
19    get_svm_detector( svm, hog_detector );
20    //set SVM
21    my_hog.setSVMDetector( hog_detector );
22
23    std::vector<Rect> positivesAll;
24    my_hog.detectMultiScale( frameGray, positivesAll, 0,
25         Size(0,0), Size(0,0), 1.1, 4);
```



Block 1  Block 2

Cells

cell histogram



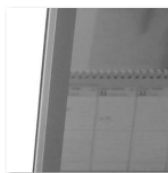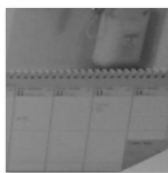https://github.com/opencv/opencv/blob/master/samples/cpp/train_HOG.cpp

# Detection step - HOG+SVM (OpenCV)

```cpp
1    int blockSize = 16;
2    int cellSize = 8;
3    int strideSize = 8;
4    int winSize = 64;
5
6    //HOGDescriptor hog;
7    HOGDescriptor my_hog(
8              cv::Size(winSize,winSize), //winSize
9              cv::Size(blockSize,blockSize), //blocksize
10             cv::Size(strideSize,strideSize), //blockStride,
11             cv::Size(cellSize,cellSize), //cellSize,
12             9, //nbins,
13             );
14
15   //SVM
16   Ptr<SVM> svm = StatModel::load<SVM>( classifierName );
17   std::vector< float > hog_detector;
18   //get the support vectors
19   get_svm_detector( svm, hog_detector );
20   //set SVM
21   my_hog.setSVMDetector( hog_detector );
22
23   std::vector<Rect> positivesAll;
24   my_hog.detectMultiScale( frameGray, positivesAll, 0,
25         Size(0,0), Size(0,0), 1.1, 4);
```
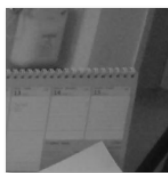


Block 1  Block 2

Cells

cell histogram

# Detection step - HOG+SVM (OpenCV)

# Object Detection/Recognition

- Haar

- HOG

- LBP

Traditional Approaches

- SIFT, SURF

KeyPoints

- CNNs

Deep Learning Approach

- Practical examples using OpenCV + Dlib (https://opencv.org/, http://dlib.net/)

# Related Works

Ahonen at al.
(2006)
1300 cit. SCOPUS

## Face Description with Local Binary Patterns: Application to Face Recognition

Timo Ahonen, *Student Member, IEEE,* Abdenour Hadid,

and Matti Pietikäinen, *Senior Member, IEEE*

Zhang at al.
(2007)

### Face Detection Based on Multi-Block LBP Representation

Lun Zhang, Rufeng Chu, Shiming Xiang, Shengcai Liao, Stan Z. Li

Center for Biometrics and Security Research & National Laboratory of Pattern Recognition
Institute of Automation, Chinese Academy of Sciences
95 Zhongguancun Donglu Beijing 100080, China

Xiaohua at al.
(2009)

## Face detection using simplified Gabor features and hierarchical regions in a cascade of classifiers

Li Xiaohua [a,b], Kin-Man Lam [b,*], Shen Lansun [c], Zhou Jiliu [a]

[a] Department of Computer Science, Sichuan University, Chengdu 610064, China
[b] Centre for Signal Processing, Department of Electronic and Information Engineering, The Hong Kong Polytechnic University, Hung Hom, Kowloon, Hong Kong
[c] Signal and Information Processing Lab., Beijing University of Technology, Beijing 100022, China

# LBP - Local Binary Patterns

- Were introduced by Ojala et al. for the texture analysis.

- The main idea behind LBP is that the local image structures (micro patterns such as lines, edges, spots, and flat areas) can be efficiently encoded by comparing every pixel with its neighboring pixels.

- Fast and cheap technique

# LBP - Local Binary Patterns

# LBP - Local Binary Patterns

- Robust to monotonic changes in illumination

# LBP - Local Binary Patterns



P = 8, R = 1.0          P = 12, R = 2.5          P = 16, R = 4.0

Ojala T, Pietikäinen M & Mäenpää T (2002) Multiresolution gray-scale and rotation invariant texture classification with Local Binary Patterns. IEEE Transactions on Pattern Analysis and Machine Intelligence 24(7):971-987

# LBP - Local Binary Patterns



$LBP_{4,1}$

$LBP_{8,1}^{u2}$

Face image divided into several blocks

LBP histogram from each block

Feature histogram

LBP histogram from the whole image

Hadid, A., Pietikainen, M., Ahonen, T.: A discriminative feature space for detecting and recognizing faces. In: Computer Vision and Pattern Recognition, 2004. CVPR 2004. Proceedings of the 2004 IEEE Computer Society Conference on. vol. 2, pp. II–797–II–804 Vol.2 (2004)

# LBP - Local Binary Patterns

Average gray-value
of Block:  7

| 6 | 8 | 8 | 8 | 12 |
| 6 | 7 | 7 | | |
| 8 | | 9 | | 11 |
| 6 | | 20 | | 19 |

Average
gray-value

Thresholding →

| 0 | 0 | 1 |
| 0 | | 1 |
| 0 | 1 | 1 |

MB-LBP: 00111100

Describing →

Zhang, L., Chu, R., Xiang, S., Liao, S., Li, S.Z.: Face detection based on multi-block lbp representation. In: Proceedings of the 2007 international conference on Advances in Biometrics. pp. 11–18. ICB'07, Springer-Verlag, Berlin, Heidelberg (2007)

# Object Detection/Recognition

- Haar

- HOG        } Traditional Approaches

- LBP

- SIFT, SURF    } KeyPoints

- CNNs        } Deep Learning Approach

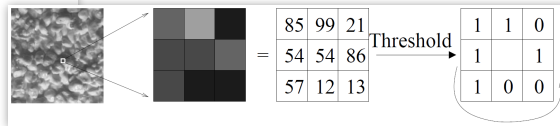- Practical examples using OpenCV + Dlib (https://opencv.org/, http://dlib.net/)

# KeyPoints

The goal is to find image KeyPoints that are invariant in the terms of scale, orientation, position, illumination, partially occlusion.

# KeyPoints – Eye Detection



template

# KeyPoints – Eye Detection

# Recognition
## Alien vs. Avenger

# Object Detection/Recognition

- Haar

- HOG    Traditional Approaches

- LBP

- SIFT, SURF    KeyPoints

- CNNs    Deep Learning Approach

- Practical examples using OpenCV + Dlib (https://opencv.org/, http://dlib.net/)

# CNNs – Main Steps (LeNet)

1. Convolution

2. Non Linearity (ReLU)

3. Pooling or Sub Sampling

4. Classification (Fully Connected Layer)



Input Image | Convolution + ReLU | Pooling | Convolution + ReLU | Pooling | Fully Connected

dog (0.01)
cat (0.04)
boat (0.94)
bird (0.02)

https://www.clarifai.com/technology

# 1. Convolution



| 1 | 1 | 1 | 0 | 0 |
|---|---|---|---|---|
| 0 | 1 | 1 | 1 | 0 |
| 0 | 0 | 1 | 1 | 1 |
| 0 | 0 | 1 | 1 | 0 |
| 0 | 1 | 1 | 0 | 0 |

Input Image

| 1 | 0 | 1 |
|---|---|---|
| 0 | 1 | 0 |
| 1 | 0 | 1 |

Mask/Filter

Convolution + ReLU   Pooling   Convolution + ReLU   Pooling   Fully Connected   Fully Connected   Output Predictions

dog (0.01)
cat (0.04)
boat (0.94)
bird (0.02)

https://ujjwalkarn.me/2016/08/11/intuitive-explanation-convnets/

# 1. Convolution

Multiply the image pixels by pixels of the filter, then sum the results

| | | |
|---|---|---|
| 1 | 0 | 1 |
| 0 | 1 | 0 |
| 1 | 0 | 1 |

Mask/Filter

| | | | | |
|---|---|---|---|---|
| 1 ×1 | 1 ×0 | 1 ×1 | 0 | 0 |
| 0 ×0 | 1 ×1 | 1 ×0 | 1 | 0 |
| 0 ×1 | 0 ×0 | 1 ×1 | 1 | 1 |
| 0 | 0 | 1 | 1 | 0 |
| 0 | 1 | 1 | 0 | 0 |

Image

| | | |
|---|---|---|
| 4 | | |
| | | |
| | | |

Convolved Feature

Convolution + ReLU    Pooling    Convolution + ReLU    Pooling    Fully Connected    Fully Connected    Output Predictions

dog (0.01)
cat (0.04)
boat (0.94)
bird (0.02)

https://ujjwalkarn.me/2016/08/11/intuitive-explanation-convnets/

# 1. Convolution

$$\text{Sharpen} - \begin{bmatrix} -1 & -1 & -1 \\ -1 & 9 & -1 \\ -1 & -1 & -1 \end{bmatrix}$$

$$\text{Blur} - \begin{bmatrix} 0 & 0.2 & 0 \\ 0.2 & 0.2 & 0.2 \\ 0 & 0.2 & 0 \end{bmatrix}$$

$$\text{Horizontal Motion Blur} - \begin{bmatrix} 0 & 0 & 0 \\ 0.2 & 0.2 & 0.2 \\ 0 & 0 & 0 \end{bmatrix}$$

$$\text{Edge detect} - \begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$$





http://dimitroff.bg/image-filtering-your-own-instagram/

# 1. Convolution



Input



Convolution
+ ReLU

Pooling

Convolution
+ ReLU

Pooling

Fully
Connected

Fully
Connected

Output Predictions

dog (0.01)
cat (0.04)
boat (0.94)
bird (0.02)

# 1. Convolution

- Before training, we have many filters/kernels
  - Filter values are randomized
- Depth of this conv. layer corresponds to the number of filters we use for the convolution operation

- The filters are learned during the training

# 2. Non Linearity (ReLU)

- ReLU is used after every Convolution operation
- The goal of this step is to replace all negative pixels by zero in the feature map



Input Feature Map — Rectified Feature Map

ReLU

Black = negative; white = positive values — Only non-negative values



Convolution + ReLU — Pooling — Convolution + ReLU — Pooling — Fully Connected — Fully Connected — Output Predictions

dog (0.01)
cat (0.04)
boat (0.94)
bird (0.02)

http://mlss.tuebingen.mpg.de/2015/slides/fergus/Fergus_1.pdf

# 3. Pooling

## ( Subsampling or downsampling )

- The goal of this step is to reduce the dimensionality of each feature map but preserve important informations

- Operations: e.g. Sum, Average, Max

# 3. Pooling

( Subsampling or downsampling )

- Common way is a pooling layer with filters of size 2x2 applied with a stride of 2

Single depth slice

# 3. Pooling

## ( Subsampling or downsampling )

- Common way is a pooling layer with filters of size 2x2 applied with a stride of 2



Max

Sum

Pooling

Only non-negative values

Rectified Feature Map



Convolution + ReLU   Pooling   Convolution + ReLU   Pooling   Fully Connected   Fully Connected   Output Predictions

dog (0.01)
cat (0.04)
boat (0.94)
bird (0.02)

# Conv. + ReLU + POOL

- Convolution layers and Pooling layers can be repeated any number of times in a single ConvNet.

# 4. Classification

- Multi Layer Perceptron
- The number of filters, filter sizes, architecture of the network etc. are fixed and do not change during training process.
- Only the values of the filter matrix and connection weights get updated.

# 4. CovNet Architectures

- **LeNet (1990s)**

- **AlexNet (2012)**

- **ZF NET (2013)**

- **GoogLeNet (2014)**

- **VGGNet (2014)**

- **ResNets (2015)**

- **DenseNet (2016)**

https://ujjwalkarn.me/2016/08/11/intuitive-explanation-convnets/

# Dlib http://dlib.net



http://dlib.net

# Recognition step CNNs (Dlib)

```cpp
using net_type = loss_multiclass_log<
                        fc<10,
                        relu<fc<84,
                        relu<fc<120,
                        max_pool<2,2,2,2,relu<con<16,5,5,1,1,
                        max_pool<2,2,2,2,relu<con<6,5,5,1,1,
                        input<matrix<unsigned char>>
                        >>>>>>>>>>>;
```

Input Image



Input Image | Convolution + ReLU | Pooling | Convolution + ReLU | Pooling | Fully Connected

dog (0.01)
cat (0.04)
boat (0.94)
bird (0.02)

http://dlib.net/dnn_introduction_ex.cpp.html

# Recognition step CNNs (Dlib)

```
using net_type = loss_multiclass_log<
                        fc<10,
                        relu<fc<84,
                        relu<fc<120,
                        max_pool<2,2,2,2,relu<con<16,5,5,1,1,
                        max_pool<2,2,2,2,relu<con<6,5,5,1,1,
                        input<matrix<unsigned char>>
                        >>>>>>>>>>>;
```

6 conv. filters
5x5 filter size
1x1 stride
+ReLU



Input Image   Convolution + ReLU   Pooling   Convolution + ReLU   Pooling   Fully Connected

dog (0.01)
cat (0.04)
boat (0.94)
bird (0.02)

http://dlib.net/dnn_introduction_ex.cpp.html

# Recognition step CNNs (Dlib)

```
using net_type = loss_multiclass_log<
                          fc<10,
                          relu<fc<84,
                          relu<fc<120,
                          max_pool<2,2,2,2,relu<con<16,5,5,1,1,
                          max_pool<2,2,2,2,relu<con<6,5,5,1,1,
                          input<matrix<unsigned char>>
                          >>>>>>>>>>>;
```

MAX POOLING
2x2 window
2x2 stride



Input Image  Convolution + ReLU  Pooling  Convolution + ReLU  Pooling  Fully Connected

dog (0.01)
cat (0.04)
boat (0.94)
bird (0.02)

http://dlib.net/dnn_introduction_ex.cpp.html

# Recognition step CNNs (Dlib)

```
using net_type = loss_multiclass_log<
                        fc<10,
                        relu<fc<84,
                        relu<fc<120,
                        max_pool<2,2,2,2,relu<con<16,5,5,1,1,
                        max_pool<2,2,2,2,relu<con<6,5,5,1,1,
                        input<matrix<unsigned char>>
                        >>>>>>>>>>;
```

16 conv. filters
5x5 filter size
1x1 stride
+ReLU



Input Image    Convolution + ReLU    Pooling    Convolution + ReLU    Pooling    Fully Connected

dog (0.01)
cat (0.04)
boat (0.94)
bird (0.02)

http://dlib.net/dnn_introduction_ex.cpp.html
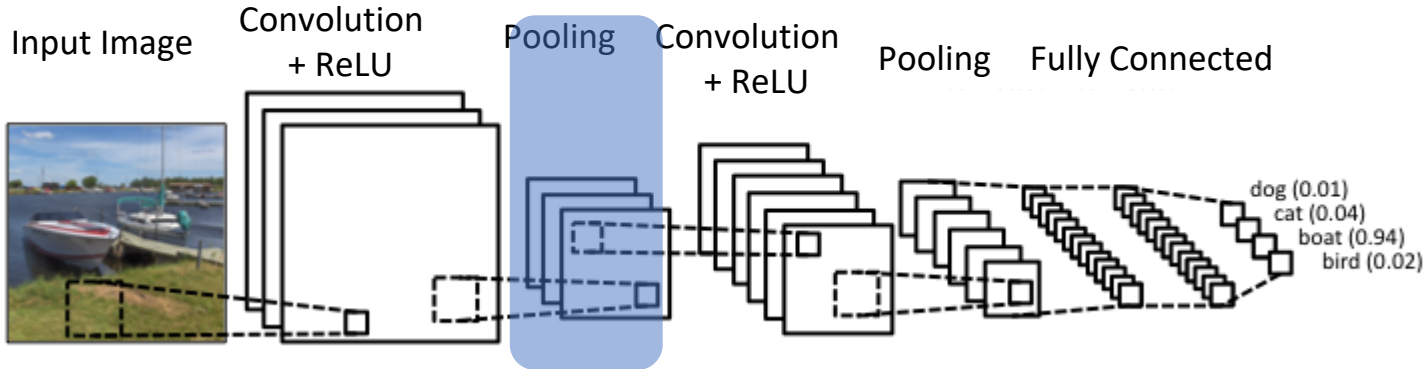
# Recognition step CNNs (Dlib)

```
using net_type = loss_multiclass_log<
                    fc<10,
                    relu<fc<84,
                    relu<fc<120,
                    max_pool<2,2,2,2,relu<con<16,5,5,1,1,
                    max_pool<2,2,2,2,relu<con<6,5,5,1,1,
                    input<matrix<unsigned char>>
                    >>>>>>>>>>>>;
```

MAX POOLING
2x2 window
2x2 stride



Input Image | Convolution + ReLU | Pooling | Convolution + ReLU | Pooling | Fully Connected

dog (0.01)
cat (0.04)
boat (0.94)
bird (0.02)

http://dlib.net/dnn_introduction_ex.cpp.html

# Recognition step CNNs (Dlib)

```
using net_type = loss_multiclass_log<
    fc<10,
    relu<fc<84,
    relu<fc<120,
    max_pool<2,2,2,2,relu<con<16,5,5,1,1,
    max_pool<2,2,2,2,relu<con<6,5,5,1,1,
    input<matrix<unsigned char>>
    >>>>>>>>>>;
```

Fully connected layer
120 neurons
84 neurons
10 outputs/classes
multiclass classification



Input Image    Convolution + ReLU    Pooling    Convolution + ReLU    Pooling    Fully Connected

dog (0.01)
cat (0.04)
boat (0.94)
bird (0.02)

http://dlib.net/dnn_introduction_ex.cpp.html

# Recognition step CNNs (Dlib)

```cpp
1    // network instance
2    net_type net;
3
4    // mini-batch stochastic gradient descent
5    //dnn_trainer<net_type> trainer(net, sgd(), {0,1}); //{0,1} - will use two GPU
6    dnn_trainer<net_type> trainer(net);
7    trainer.set_learning_rate(0.01);
8    trainer.set_min_learning_rate(0.0001);
9    trainer.set_mini_batch_size(160);
10   trainer.set_iterations_without_progress_threshold(500);
11   trainer.set_max_num_epochs(100);
12   trainer.be_verbose();
13   //train
14   trainer.train(train_images, train_labels);
15   // save
16   serialize("LeNet.dat") << net;
```

http://dlib.net/dnn_introduction_ex.cpp.html

# Recognition step CNNs (Dlib + OpenCV)

```cpp
1  //Load image using OpenCV
2  Mat frame;
3  frame = imread( "my_img.png", 1 );
4  cvtColor( frame, frame, COLOR_BGR2GRAY );
5  medianBlur(frame, frame, 5);
6
7   //OpenCV Mat to Dlib
8  cv_image<unsigned char> cimg(frame);
9  matrix<unsigned char> dlibFrame = dlib::mat(cimg);
10
11 //prediction using CNN
12 unsigned long predict_label = net(frame);
```

http://dlib.net/dnn_introduction_ex.cpp.html

# Recognition step CNNs (dlib)

# CNNs (Dlib)

**NVIDIA 1080ti - 39 frames per second, 928x478**

# Thank you for your attention

## http://mrl.cs.vsb.cz