# Histograms of Oriented Gradients for Human Detection (HOG)

N. Dalal and B. Triggs

CVPR 2005

# 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

## Abstract

*We study the question of feature sets for robust visual object recognition, adopting linear SVM based human detection as a test case. After reviewing existing edge and gradient based descriptors, we show experimentally that grids of Histograms of Oriented Gradient (HOG) descriptors significantly outperform existing feature sets for human detection. We study the influence of each stage of the computation on performance, concluding that fine-scale gradients, fine orientation binning, relatively coarse spatial binning, and high-quality local contrast normalization in overlapping descriptor blocks are all important for good results. The new approach gives near-perfect separation on the original MIT pedestrian database, so we introduce a more challenging dataset containing over 1800 annotated human images with a large range of pose variations and backgrounds.*

## 1 Introduction

Detecting humans in images is a challenging task owing to their variable appearance and the wide range of poses that they can adopt. The first need is a robust feature set that allows the human form to be discriminated cleanly, even in cluttered backgrounds under difficult illumination. We study the issue of feature sets for human detection, showing that locally normalized Histogram of Oriented Gradient (HOG) descriptors provide excellent performance relative to other existing feature sets including wavelets [17,22]. The proposed descriptors are reminiscent of edge orientation histograms [4,5], SIFT descriptors [12] and shape contexts [1], but they are computed on a dense grid of uniformly spaced cells and they use overlapping local contrast normalizations for improved performance. We make a detailed study of the effects of various implementation choices on detector performance, taking "pedestrian detection" (the detection of mostly visible people in more or less upright poses) as a test case. For simplicity and speed, we use linear SVM as a baseline classifier throughout the study. The new detectors give essentially perfect results on the MIT pedestrian test set [18,17], so we have created a more challenging set containing over 1800 pedestrian images with a large range of poses and backgrounds. Ongoing work suggests that our feature set performs equally well for other shape-based object classes.

We briefly discuss previous work on human detection in §2, give an overview of our method §3, describe our data sets in §4 and give a detailed description and experimental evaluation of each stage of the process in §5–6. The main conclusions are summarized in §7.

## 2 Previous Work

There is an extensive literature on object detection, but here we mention just a few relevant papers on human detection [18,17,22,16,20]. See [6] for a survey. Papageorgiou et al [18] describe a pedestrian detector based on a polynomial SVM using rectified Haar wavelets as input descriptors, with a parts (subwindow) based variant in [17]. Depoortere et al give an optimized version of this [2]. Gavrila & Philomen [8] take a more direct approach, extracting edge images and matching them to a set of learned exemplars using chamfer distance. This has been used in a practical real-time pedestrian detection system [7]. Viola et al [22] build an efficient moving person detector, using AdaBoost to train a chain of progressively more complex region rejection rules based on Haar-like wavelets and space-time differences. Ronfard et al [19] build an articulated body detector by incorporating SVM based limb classifiers over 1st and 2nd order Gaussian filters in a dynamic programming framework similar to those of Felzenszwalb & Huttenlocher [3] and Ioffe & Forsyth [9]. Mikolajczyk et al [16] use combinations of orientation-position histograms with binary-thresholded gradient magnitudes to build a parts based method containing detectors for faces, heads, and front and side profiles of upper and lower body parts. In contrast, our detector uses a simpler architecture with a single detection window, but appears to give significantly higher performance on pedestrian images.

## 3 Overview of the Method

This section gives an overview of our feature extraction chain, which is summarized in fig. 1. Implementation details are postponed until §6. The method is based on evaluating well-normalized local histograms of image gradient orientations in a dense grid. Similar features have seen increasing use over the past decade [4,5,12,15]. The basic idea is that local object appearance and shape can often be characterized rather well by the distribution of local intensity gradients or

3

# HOG Steps

- HOG feature extraction
  - Compute centered horizontal and vertical gradients with no smoothing
  - Compute gradient orientation and magnitudes
    - For color image, pick the color channel with the highest gradient magnitude for each pixel.

  - For a 64x128 image,
  - Divide the image into 16x16 blocks of 50% overlap.
    - 7x15=105 blocks in total
  - Each block should consist of 2x2 cells with size 8x8.
  - Quantize the gradient orientation into 9 bins
    - The vote is the gradient magnitude
    - Interpolate votes between neighboring bin center.
    - The vote can also be weighted with Gaussian to downweight the pixels near the edges of the block.
  - Concatenate histograms (Feature dimension: 105x4x9 = 3,780)

# Computing Gradients

- **Centered:** $f'(x) = \lim_{h \to 0} \dfrac{f(x+h) - f(x-h)}{2h}$
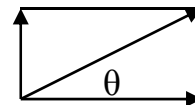
- **Filter masks in x and y directions**

  - Centered:

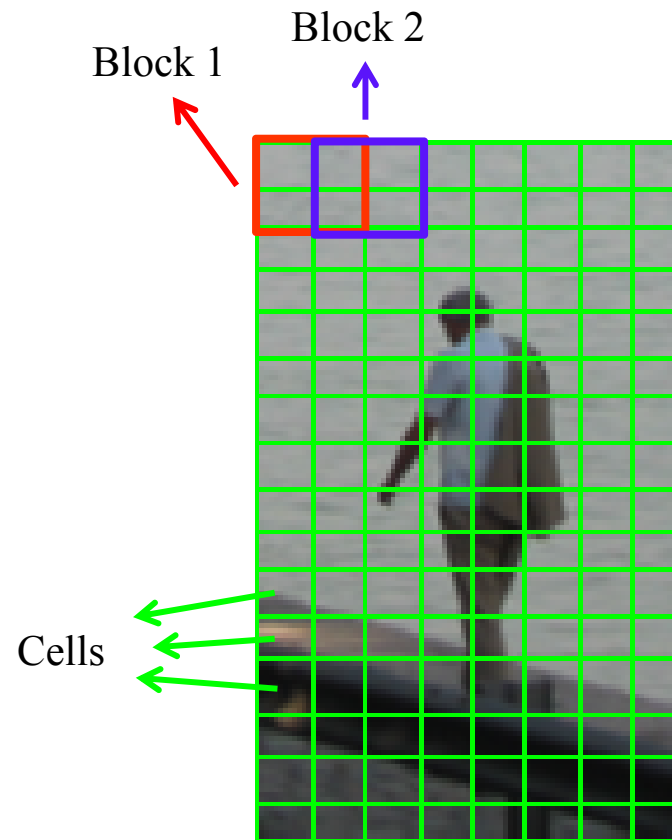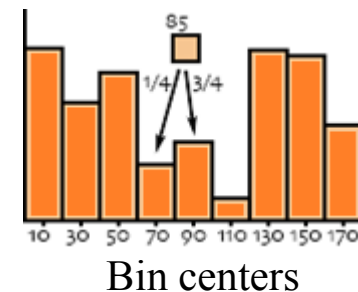    | -1 | 0 | 1 |
    |----|---|---|

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

- **Gradient**

  - Magnitude: $s = \sqrt{s_x^2 + s_y^2}$

  - Orientation: $\theta = \arctan(\dfrac{s_y}{s_x})$

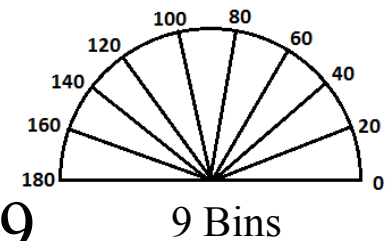# Blocks, Cells

- 16x16 blocks of 50% overlap.
  - 7x15=105 blocks in total

- Each block should consist of 2x2 cells with size 8x8.

Block 1
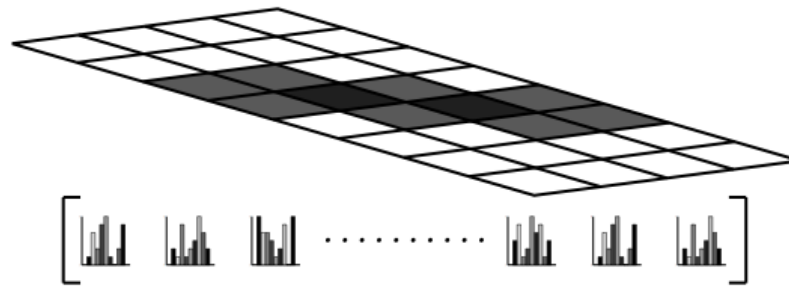
Block 2



Cells

# Votes

- Each block consists of 2x2 cells with size 8x8

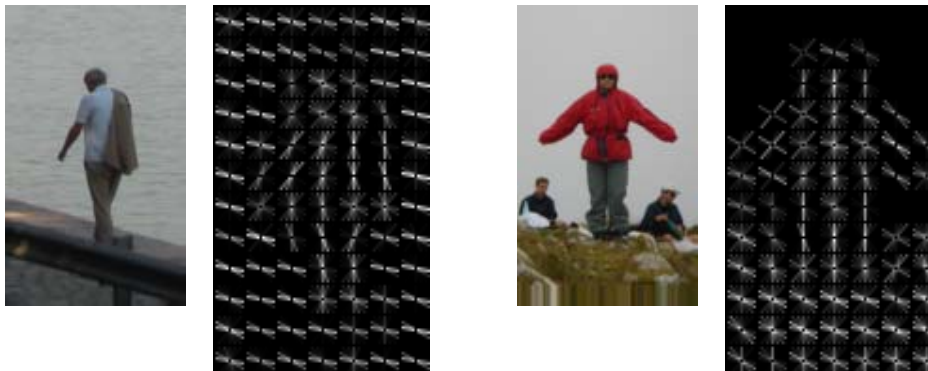- Quantize the gradient orientation into 9 bins (0-180)

  - The vote is the gradient magnitude

  - Interpolate votes linearly between neighboring bin centers.

    - Example: if $\theta$=85 degrees.

    - Distance to the bin center Bin 70 and Bin 90 are 15 and 5 degrees, respectively.

    - Hence, ratios are 5/20=1/4, 15/20=3/4.
    .

9 Bins

Bin centers

# Final Feature Vector

- Concatenate histograms
  - Make it a 1D vector of length 3780.



- Visualization

# Results

Navneet Dalal and Bill Triggs "Histograms of Oriented Gradients for Human Detection" CVPR05

# SIFT Vs HOG

**SIFT**

- 128 dimensional vector
- 16 by 16 window
- 4x4 sub-window (16 total)
- 8 bin histogram

**HOG**

- 3,780 dimensional vector
- 64 by 128 window
- 16 by 16 blocks with overlap
- Each block consists of 2 by 2 cells each of 8 by 8
- Overlapping
- 9 bin histogram