

Project Presentations

- Feedback on proposals is up
- 67 individuals/groups
- 5 min Presentation + 3 min Q & A
- $536 \text{ min} \approx 9 \text{ hours}$ AND **no** factor of safety
- Initially planned dates: Nov 27 (50') + Nov 29 (100') + Dec 4 (50') = 200'
- Only option: extra session on Dec 6 (100')
 - That will give us 300' total = 4' per presentation, and no Q & A
- Will randomly assign slots (You only need to show up for your own session)
 - No walking in and out please. Come for the whole session.
 - You are of course, encouraged to attend the other 3 sessions too
- The presentation would be created as a MarkUs deliverable. Upload as a PDF latest by Noon the day your presentation is scheduled.
- Connecting/disconnecting individual laptops would waste a lot of time

Project Presentations

- Or maybe
 - Instead of in-class presentations, you each make a video
 - Due on 3rd for everyone (instead of Nov 27, 29, Dec 4, 6, randomly assigned)
 - Submit your slides (PDF) and a link to your video on MarkUs
 - You can share as an unlisted (not private) YouTube video, or however you like

Let's vote

Object Detection

Sliding Windows

Type of Approaches

Different approaches tackle detection differently. They can roughly be categorized into three main types:

- Find **interest points**, followed by Hough voting
- **Sliding windows**: “slide” a box around image and classify each image crop inside a box (contains object or not?) ← **Let's look at a few methods for this**
- Generate **region (object) proposals**, and classify each region

Sliding Window Approaches

There are many... We will look at two in more detail:

- Dalal and Triggs (2005): HOG (Person) Detector (24,400+ citations)
- Felzenswalb et al. (2010): Deformable Part-based Model (7,400+ citations)

The last detector (DPM) is an extension of Dalal & Triggs. Another famous detector Viola Jones – we will not cover, but read the paper):

- Viola and Jones (2001): (Face) Detector (17,900+ citations)

Sliding Window Approaches

There are many... We will look at two in more detail:

- Dalal and Triggs (2005): HOG (Person) Detector → This first
- Felzenswalb et al. (2010): Deformable Part-based Model

The HOG Detector

N. Dalal and B. Triggs

Histograms of oriented gradients for human detection

CVPR, 2005

Paper: <http://lear.inrialpes.fr/people/triggs/pubs/Dalal-cvpr05.pdf>

The HOG Detector

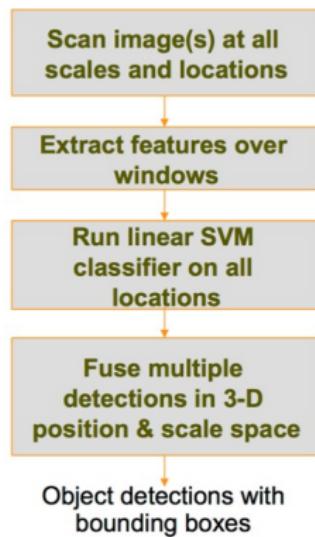
- We want to find all people in this image. Preferably our detections should not include trees, lamp posts and umbrellas.



The HOG Detector

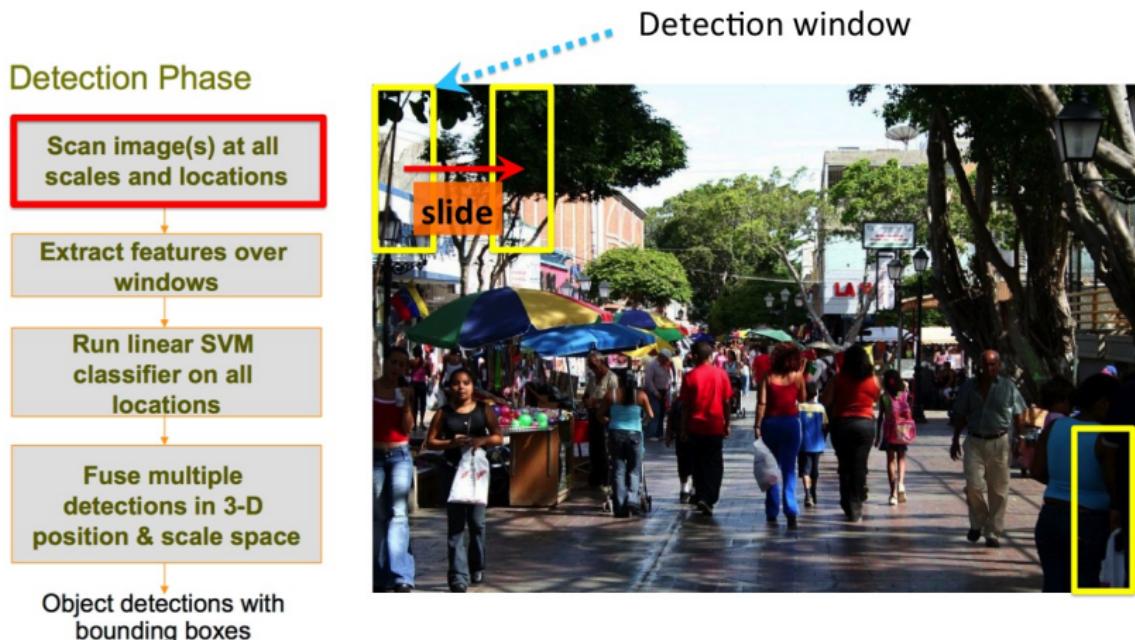
- Sliding window detectors find objects in 4 very simple steps: **(1.)** inspect every window, **(2.)** extract features in window, **(3.)** classify & accept wind. if score above threshold, **(4.)** clean-up the mess (called post-processing)

Detection Phase



The HOG Detector – Sliding the Window

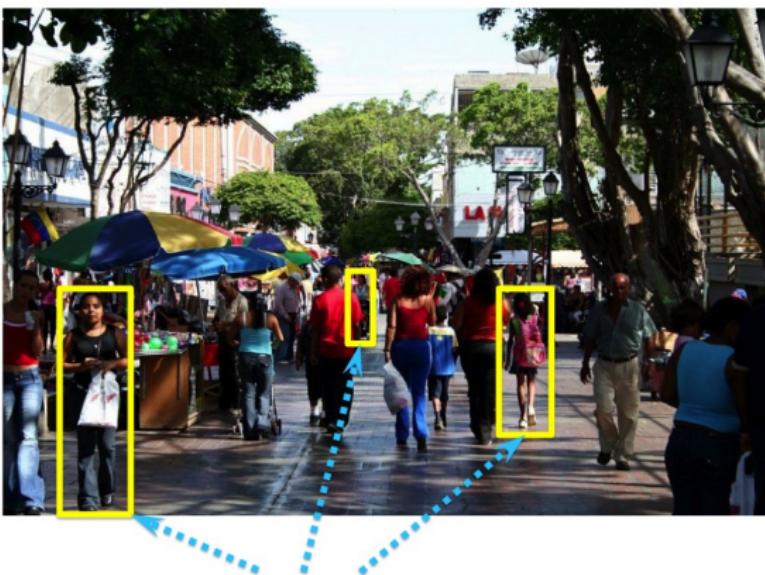
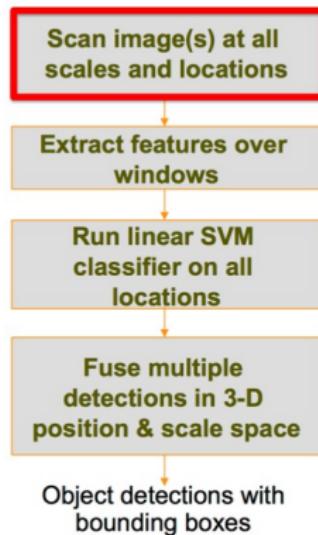
- First step: inspect every window. Typically the size of window is **fixed**.



The HOG Detector – Sliding the Window

- Since window size is fixed, how can we find people at different sizes?

Detection Phase

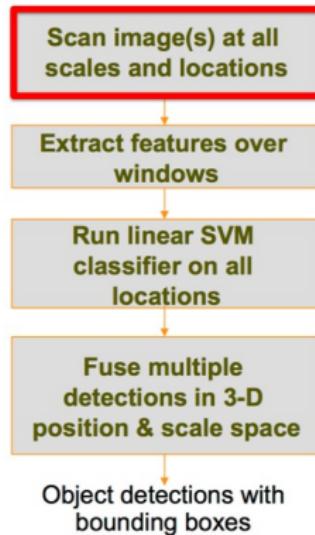


Objects can be of very different sizes (scales), even in the same image. How do we deal with that?

The HOG Detector – Sliding the Window

- Shrink (down-scale) the image and slide again

Detection Phase

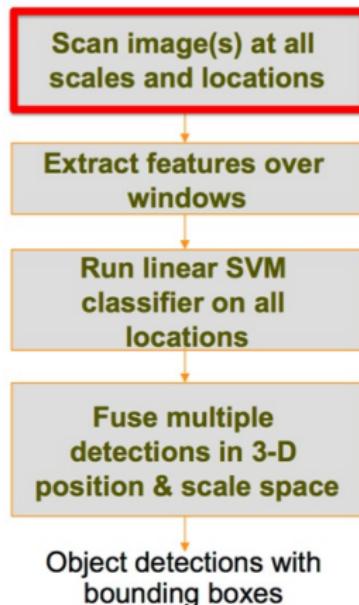


Scale-down the image, and slide the window again (the size of the window is always the same)

The HOG Detector – Sliding the Window

- Keep shrinking and sliding

Detection Phase

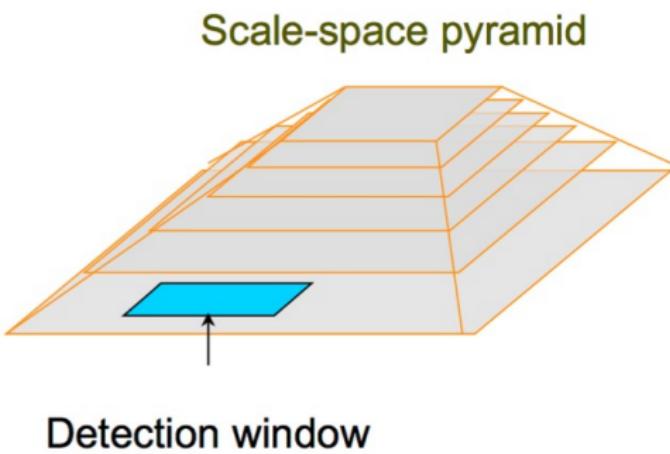
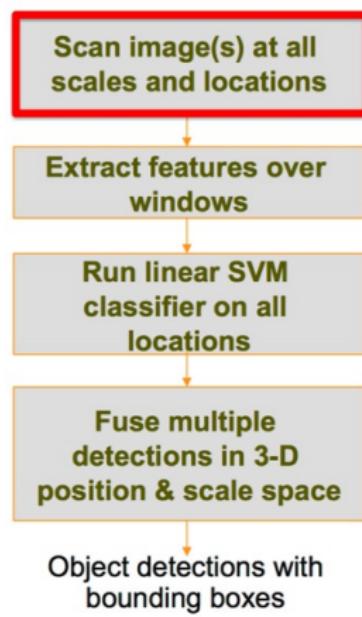


And again...

The HOG Detector – Sliding the Window

- In fact, do a full image pyramid, and slide your detector at each scale. Make sure the scale differences across levels are small (do lots of re-scaled images)

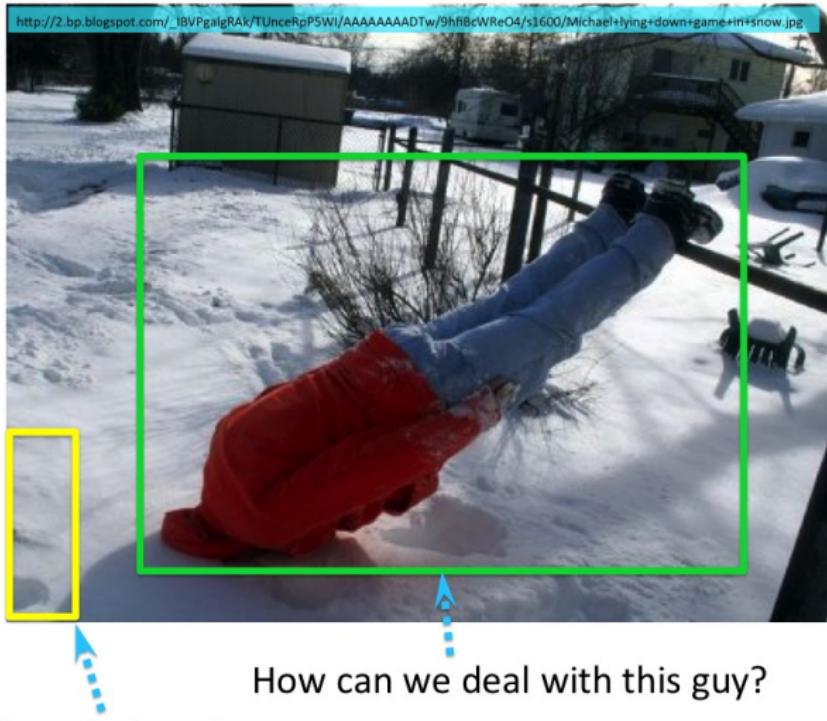
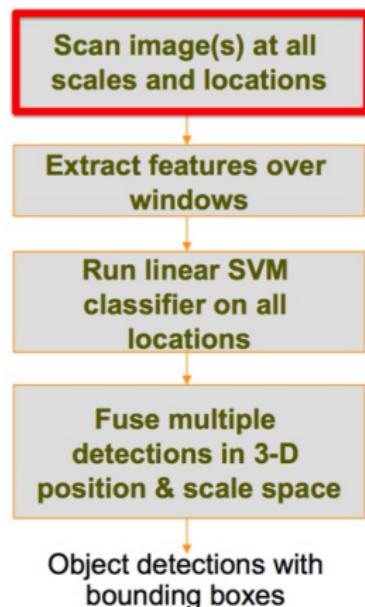
Detection Phase



The HOG Detector – Sliding the Window?

- What if the object is in a weird pose (window is of different aspect ratio)?

Detection Phase



Our window size

How can we deal with this guy?

The HOG Detector – Limitations

- Stop thinking too hard. In 2005 people were only in upright position.
- We will re-visit this question a little later (when we talk about DPM)

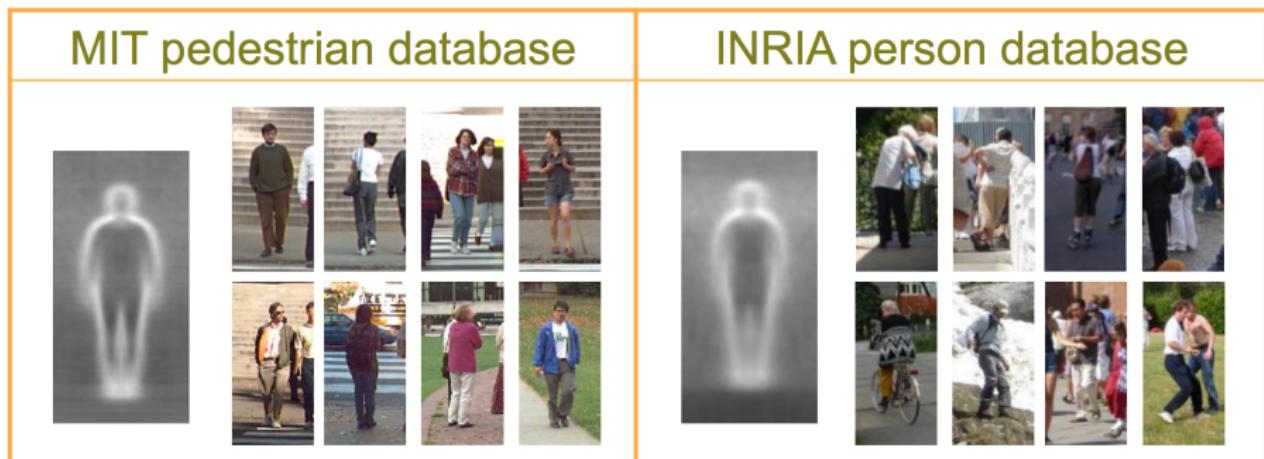
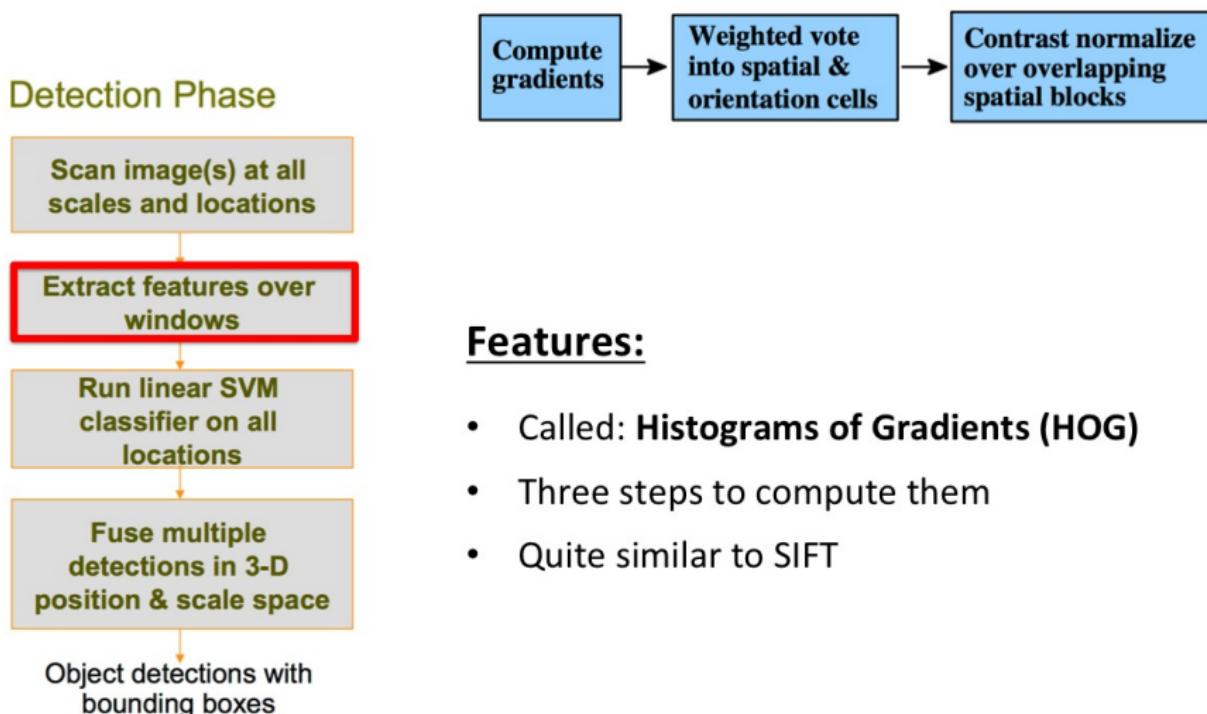


Figure: Main pedestrian detection datasets prior to PASCAL VOC.

The HOG Detector – Features (HOG)

- Famous feature descriptor called HOG that replaced SIFT (at least for object detection). There are three steps to compute it.



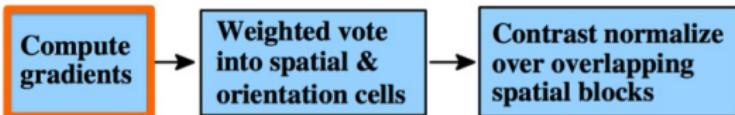
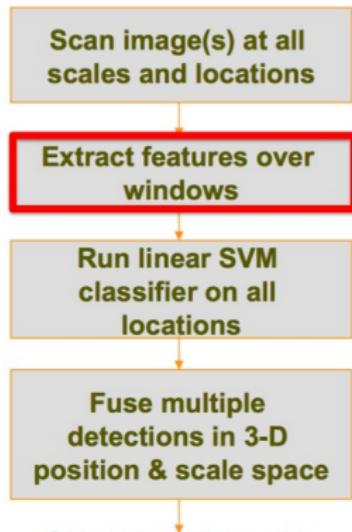
Features:

- Called: **Histograms of Gradients (HOG)**
- Three steps to compute them
- Quite similar to SIFT

The HOG Detector – Features (HOG)

- First compute gradients

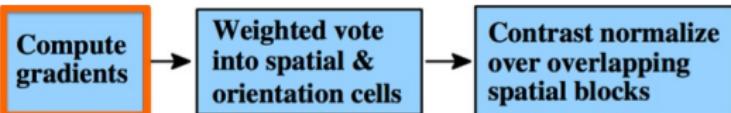
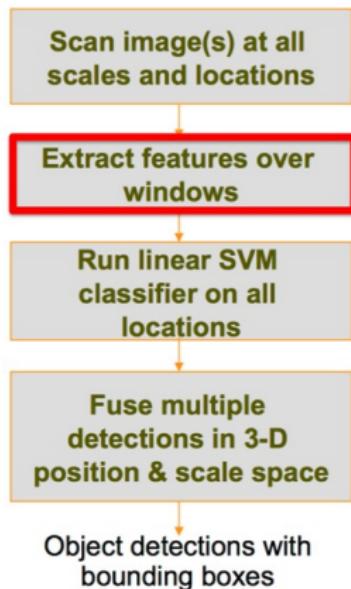
Detection Phase



The HOG Detector – Features (HOG)

- There are many ways how to compute the gradients. The HOG detector guys tried a lot of them and picked the best one.

Detection Phase



Mask Type	1D centered	1D uncentered	1D cubic-corrected	2x2 diagonal	3x3 Sobel
Operator	$[-1, 0, 1]$	$[-1, 1]$	$[1, -8, 0, 8, -1]$	$\begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}$	$\begin{bmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{bmatrix}$
Miss rate at 10^{-4} FPPW	11%	12.5%	12%	12.5%	14%

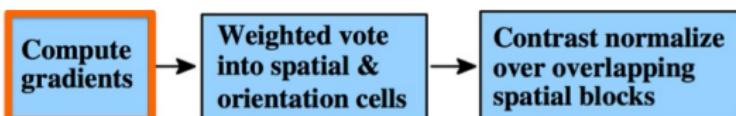
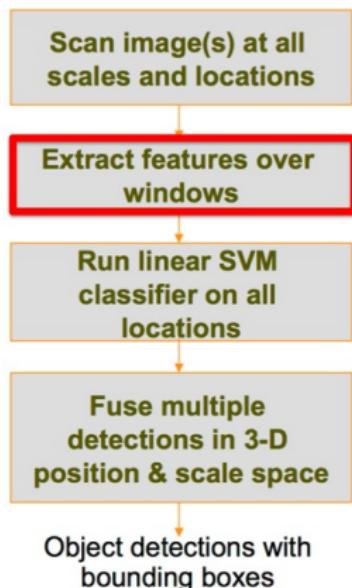
(Miss rate: smaller is better)

This gradient filter gives the best performance

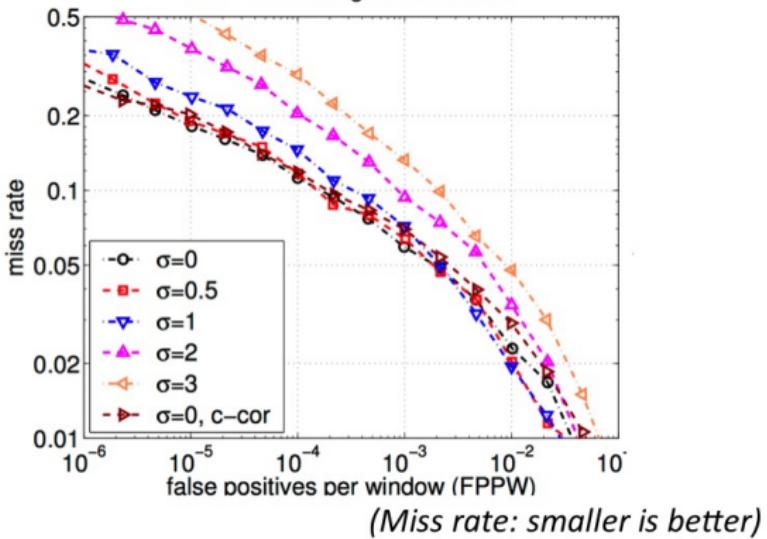
The HOG Detector – Features (HOG)

- One can also smooth image before computing the gradients. The HOG detector guys tested that as well.

Detection Phase



DET – effect of gradient scale σ

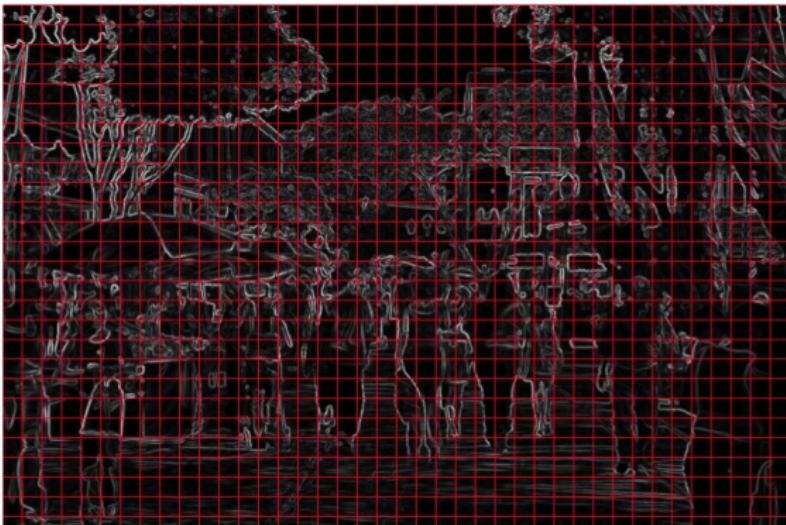
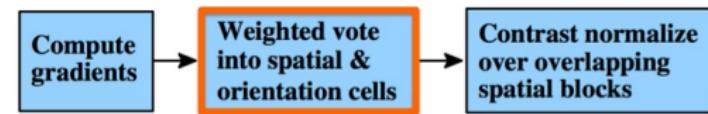
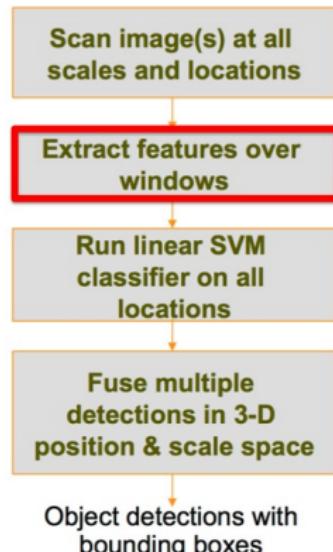


No Gaussian smoothing gives the best performance

The HOG Detector – Features (HOG)

- Divide the image into **cells** of 8×8 pixels.

Detection Phase

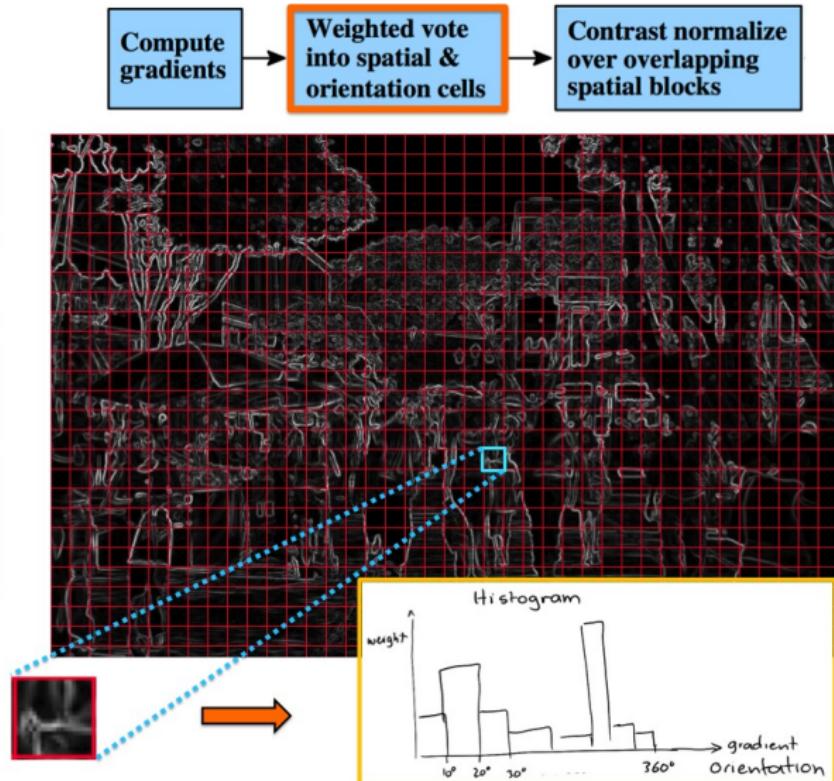
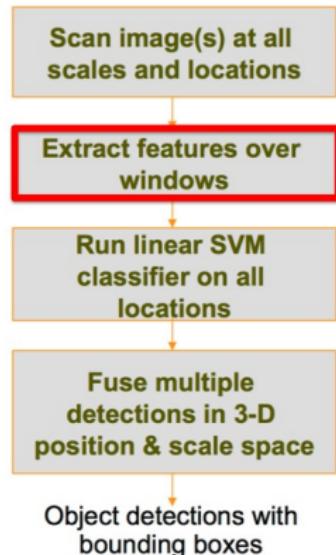


Divide the gradient image into non-overlapping **cells**.
Each cell is typically 8×8 pixels.

The HOG Detector – Features (HOG)

- Compute a histogram of orientations in each cell (similar to SIFT)

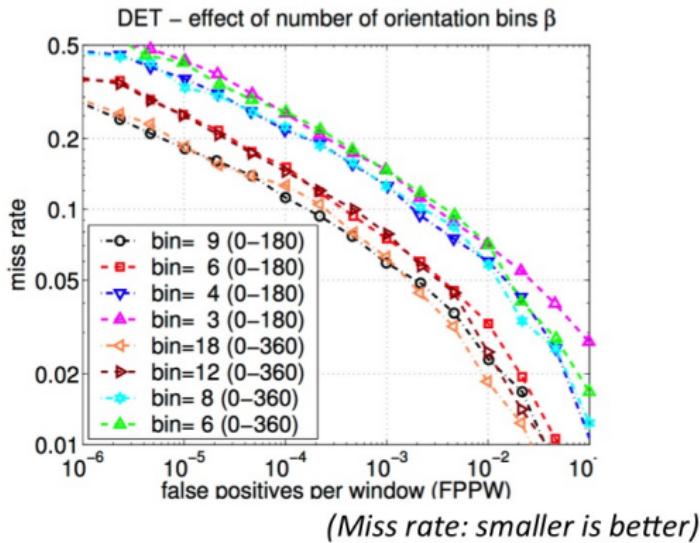
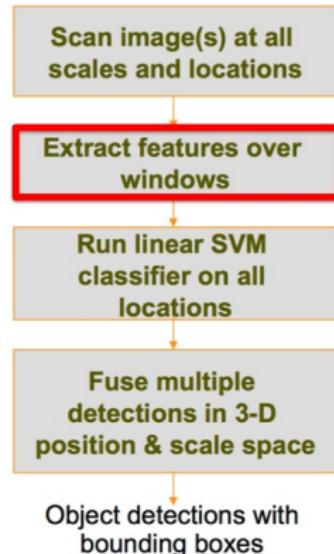
Detection Phase



The HOG Detector – Features (HOG)

- Again, check how many bins is best to use. Turns out: 9 with orient 0-180.

Detection Phase

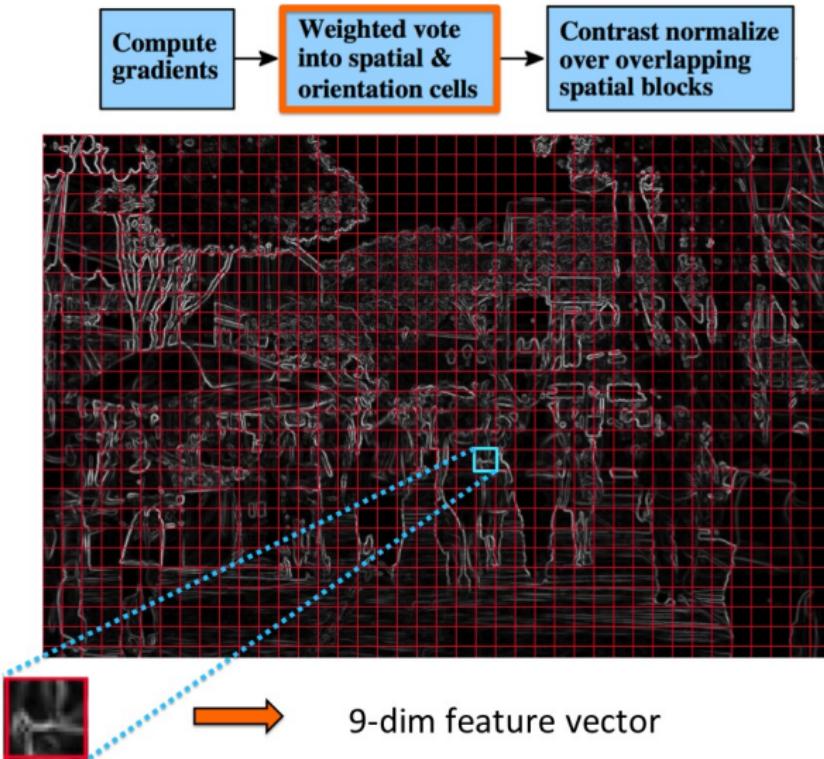
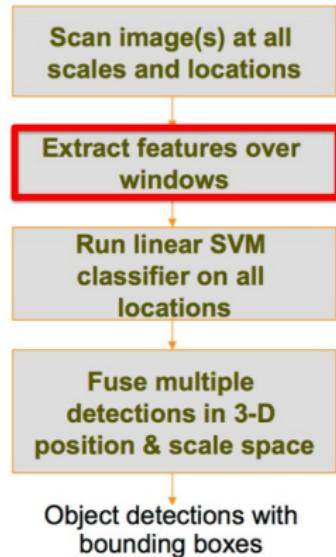


9 bins (unsigned orient) is best

The HOG Detector – Features (HOG)

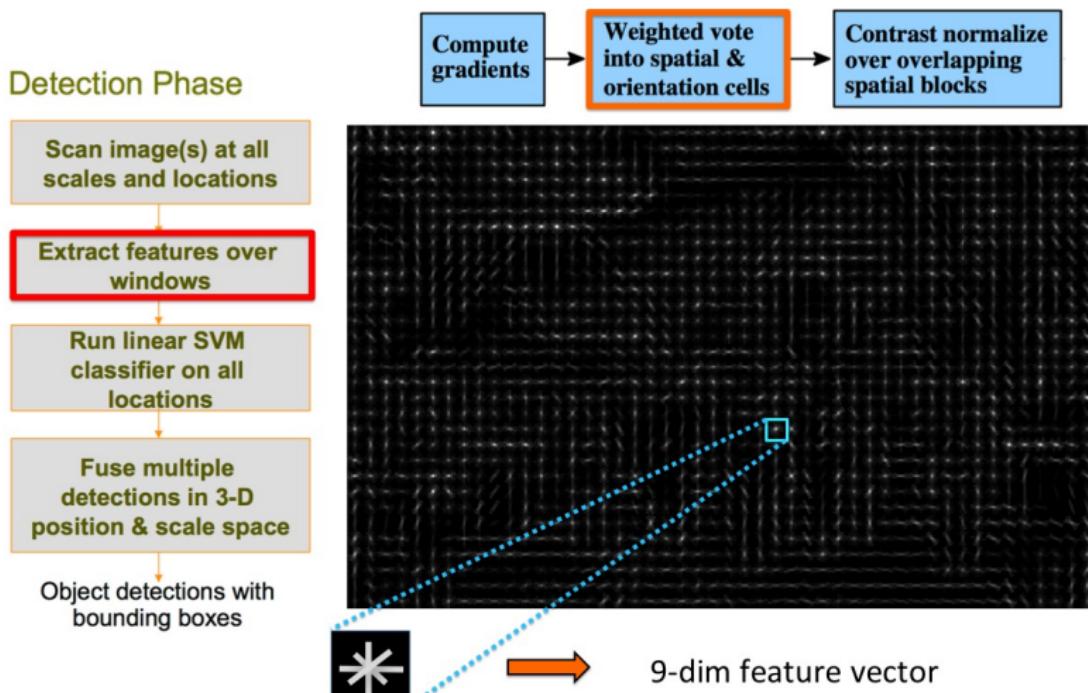
- So each cell now has a 9-dimensional feature vector

Detection Phase



The HOG Detector – Features (HOG)

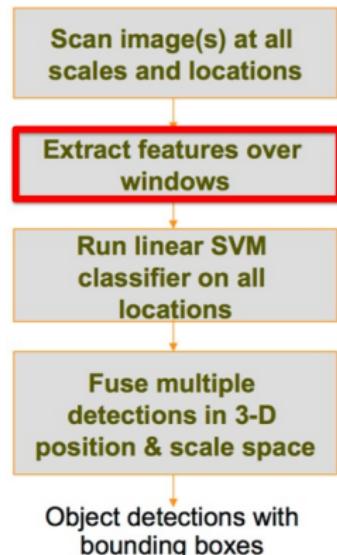
- In literature you will see this kind of **visualization** for HOG. In each cell people plot all the orientations that are present in the cell. Do not confuse this visualization with the actual feature (composed of 9 elements).



The HOG Detector – Features (HOG)

- We're not finished. We now take **blocks**, where each block has 2×2 cells.

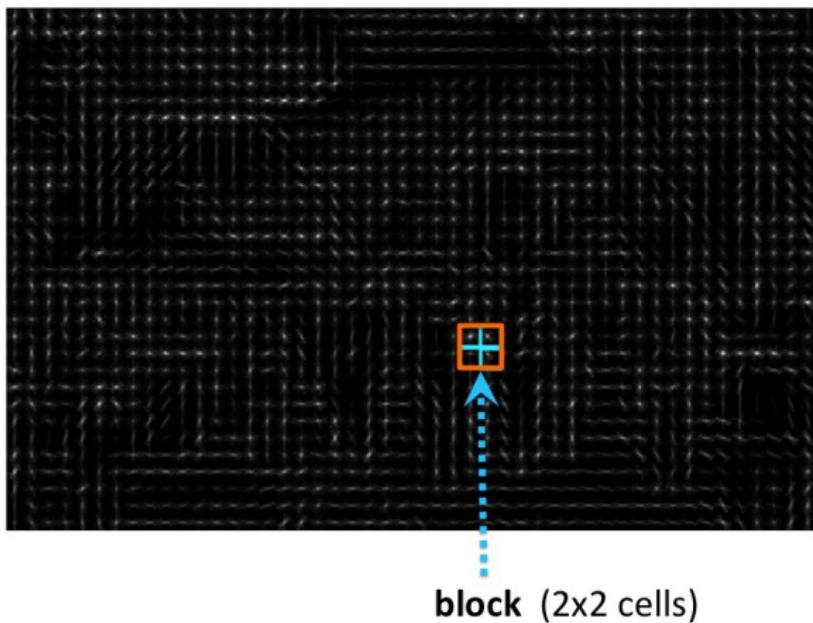
Detection Phase



Compute gradients

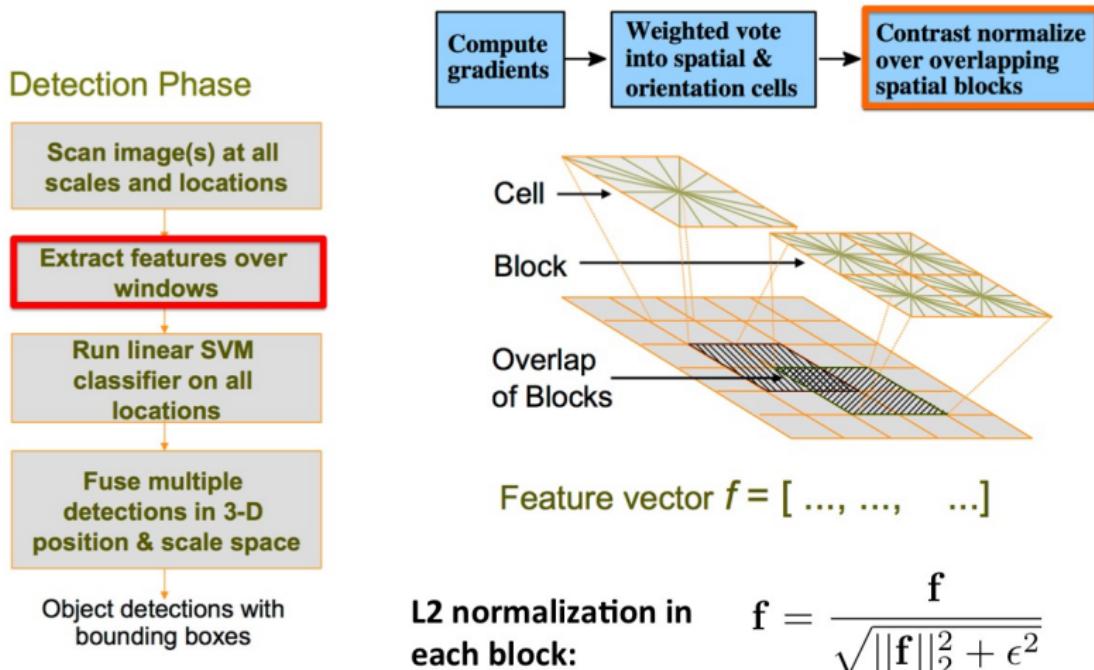
Weighted vote into spatial & orientation cells

Contrast normalize over overlapping spatial blocks



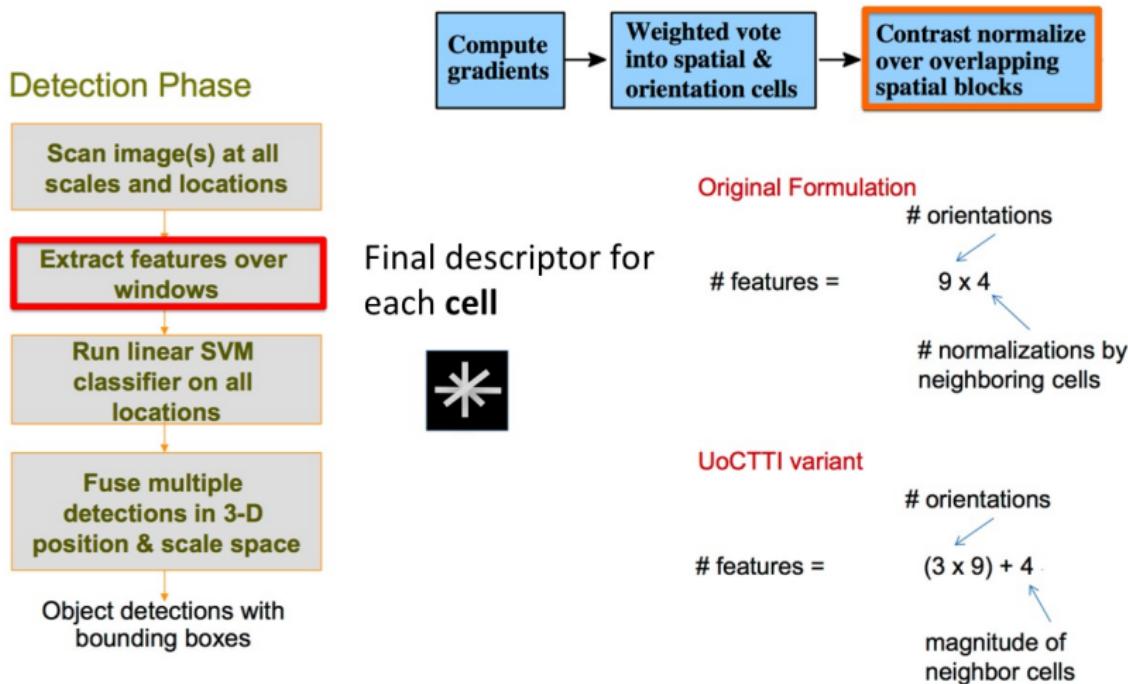
The HOG Detector – Features (HOG)

- We normalize each feature vector, such that each block has unit norm. This step doesn't change the dimension of the feature, just the strength. Why are we doing this?



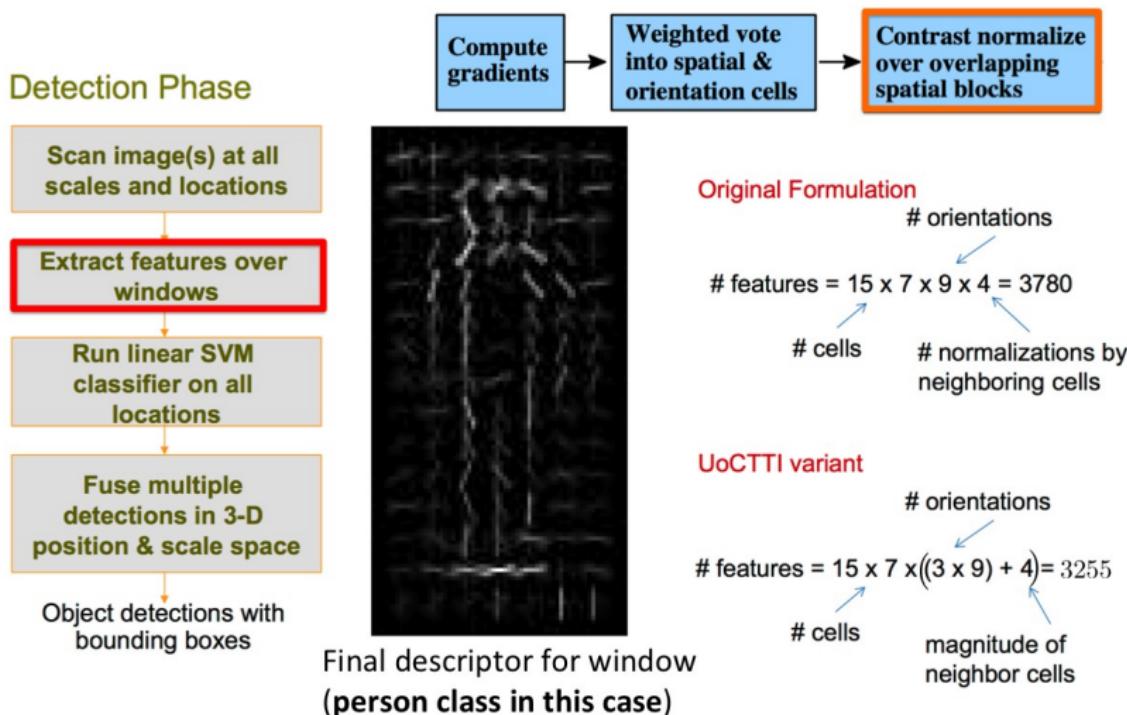
The HOG Detector – Features (HOG)

- Since each cell is in 4 blocks, we have 4 different normalizations, and we make each one into separate features.



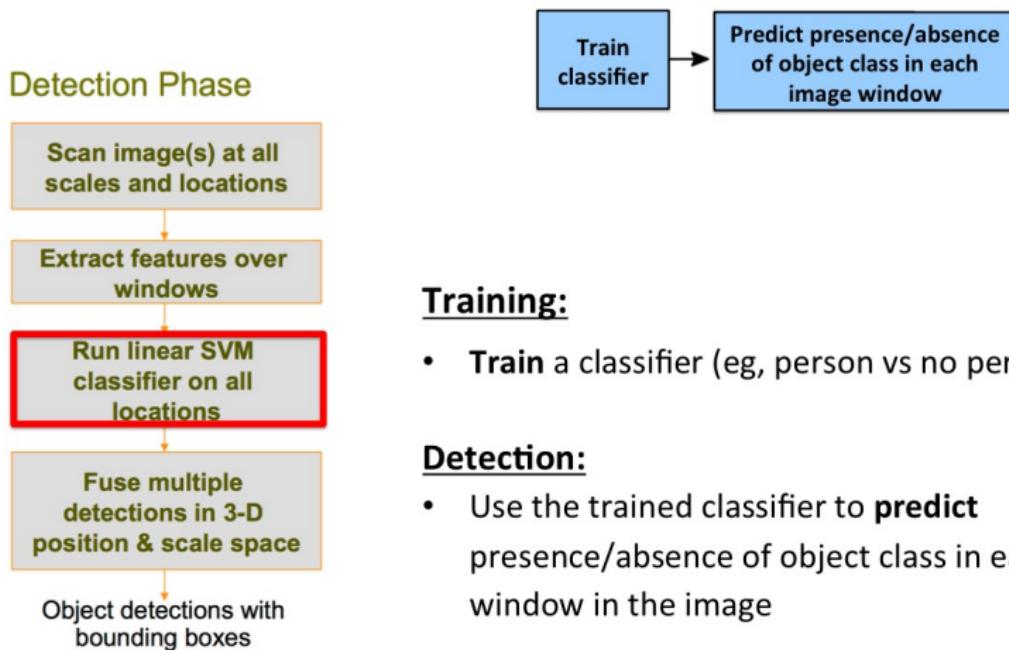
The HOG Detector – Features (HOG)

- For person class, window is 15×7 HOG cells (what's the size in pixels?)
- We vectorize the feature matrix in each window.



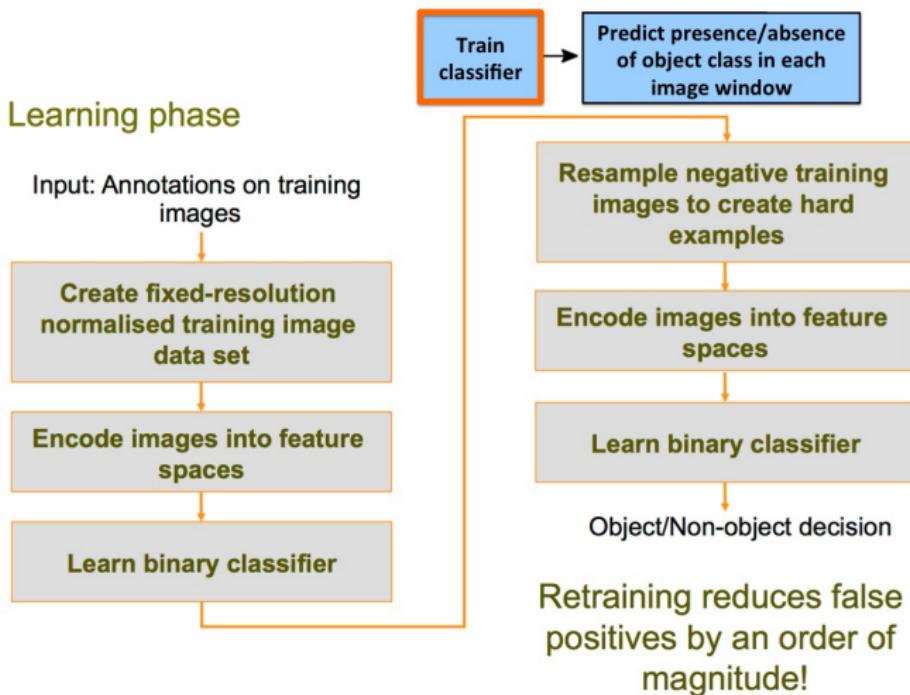
The HOG Detector – Classification

- Features done, we are ready for classification. We first need to **train** our classifier, and only after we can do detection (prediction).



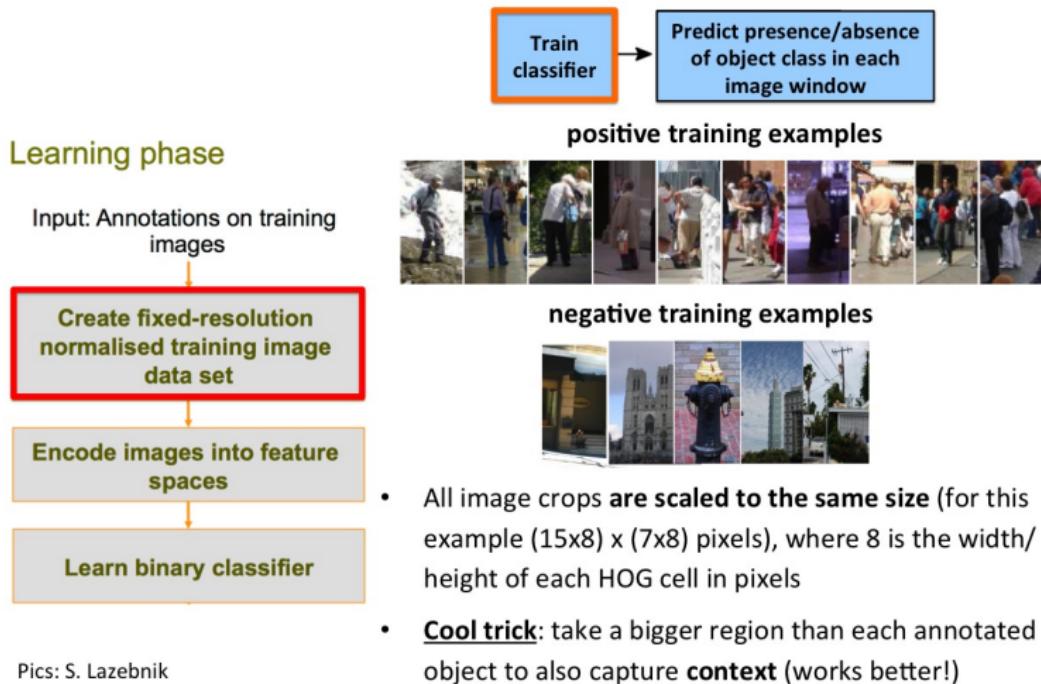
The HOG Detector – Training

- Several simple steps. Plus a few useful additional tricks (remember, some hacking is part of a Vision Researcher's life).



The HOG Detector – Training

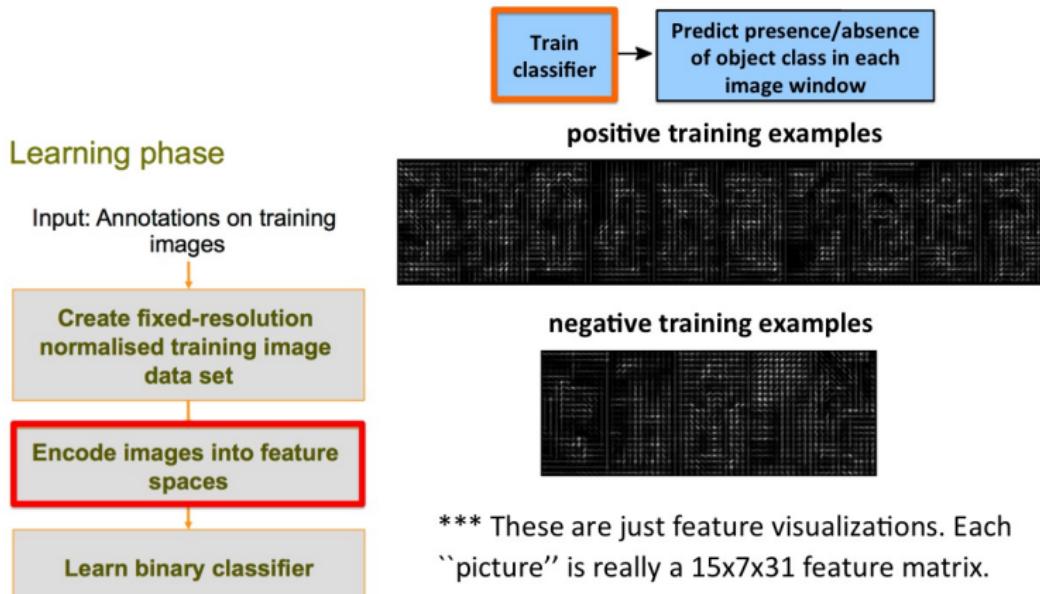
- Take a dataset with annotations. If nothing exists, collect and label yourself.



Pics: S. Lazebnik

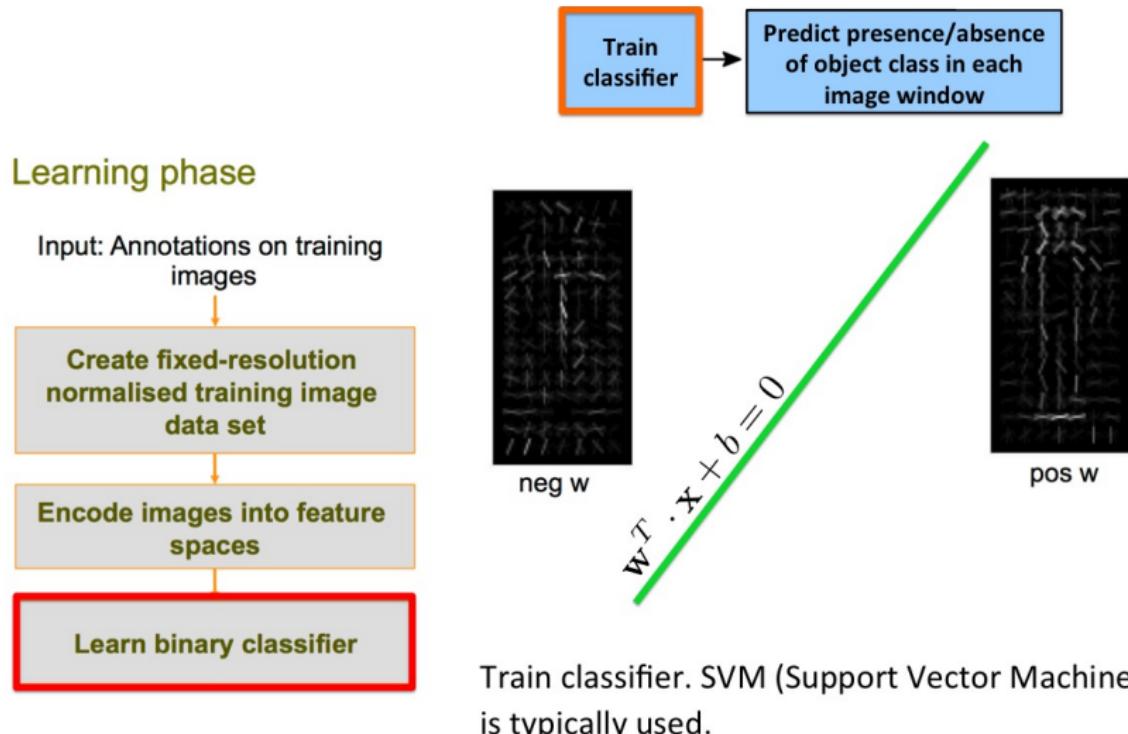
The HOG Detector – Training

- Scale positive and negative examples to the size of detection window.
Compute HOG.



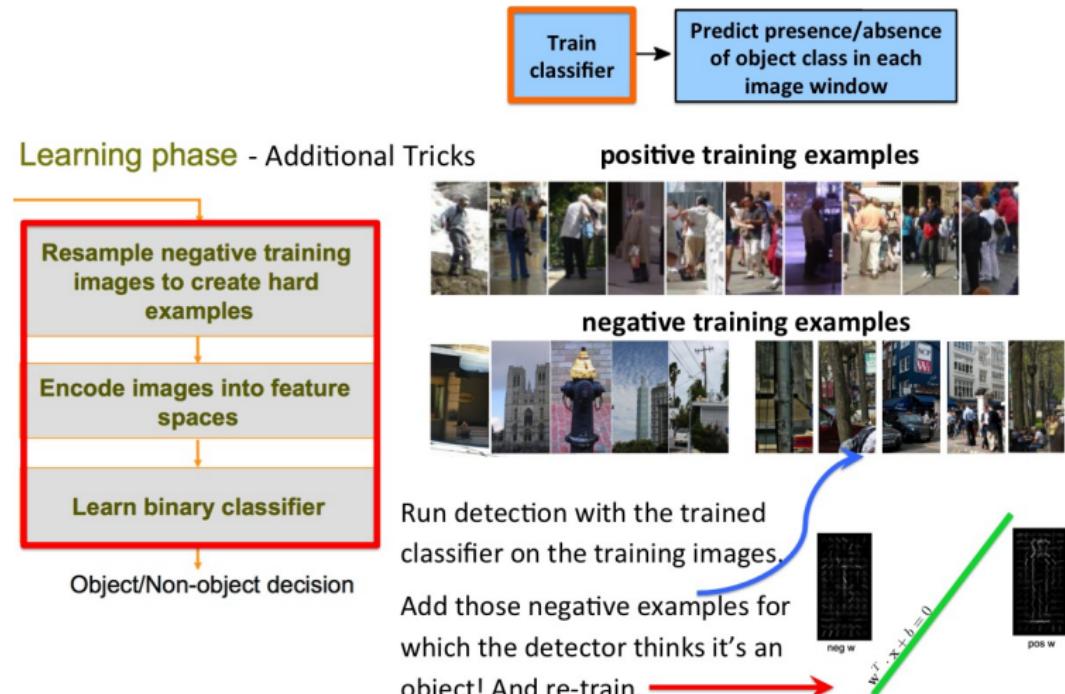
The HOG Detector – Training

- Train a classifier (with e.g. LibSVM).



The HOG Detector – Training

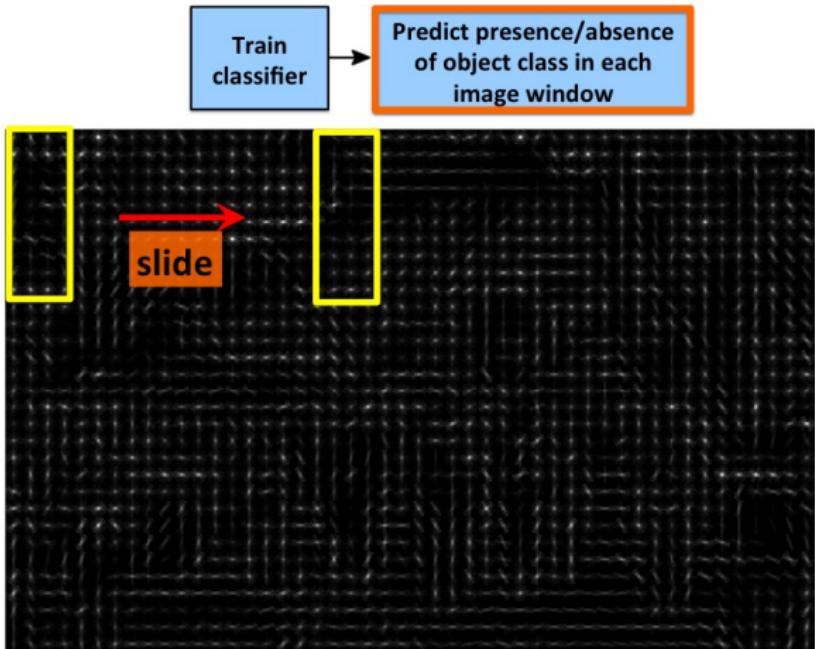
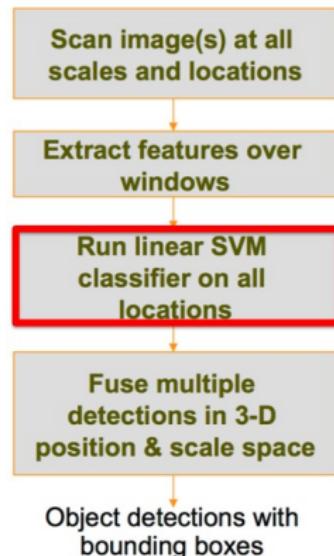
- Additional tricks: **Bootstrapping**. A fancy name for running your classifier on **training** images (with full detection pipeline), and finding mis-classified windows. Add those to training examples, and re-train classifier.



The HOG Detector – Detection

- Take a window, crop out a feature matrix, vectorize and classify

Detection Phase

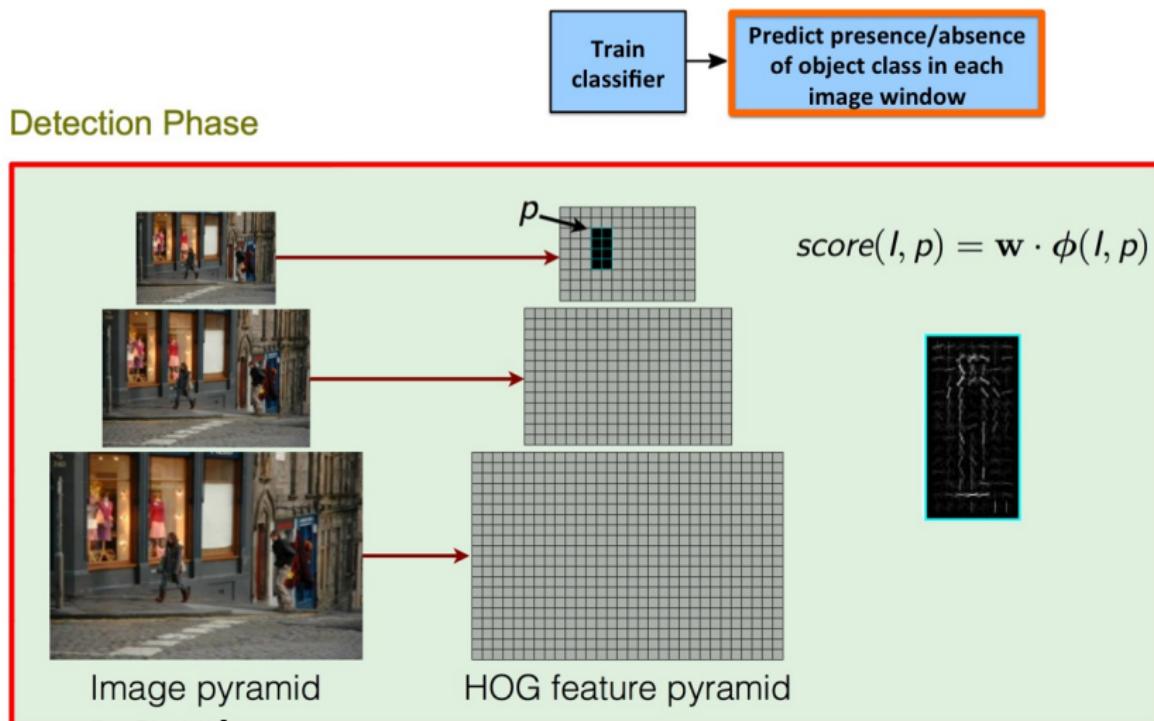


Crop out a feature $\mathbf{x} = \mathbf{f}(\cdot)$ for each window

Compute: $\text{score} = \mathbf{w}^T \cdot \mathbf{x} + b$ (higher better)

The HOG Detector – Detection

- Computing the score $\mathbf{w}^T \cdot \mathbf{x} + b$ in every location is the same as performing **cross-correlation with template \mathbf{w}** (and add b to result).

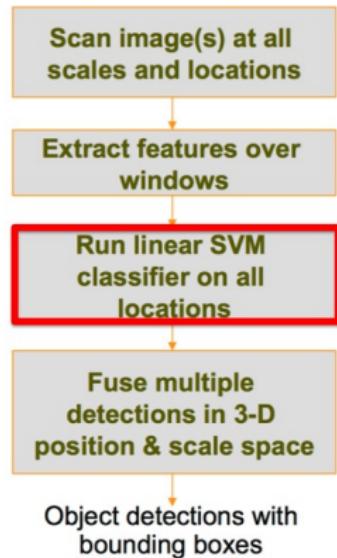


[Pic from: R. Girshik]

The HOG Detector – Training

- Threshold the scores (e.g., score > -1)

Detection Phase

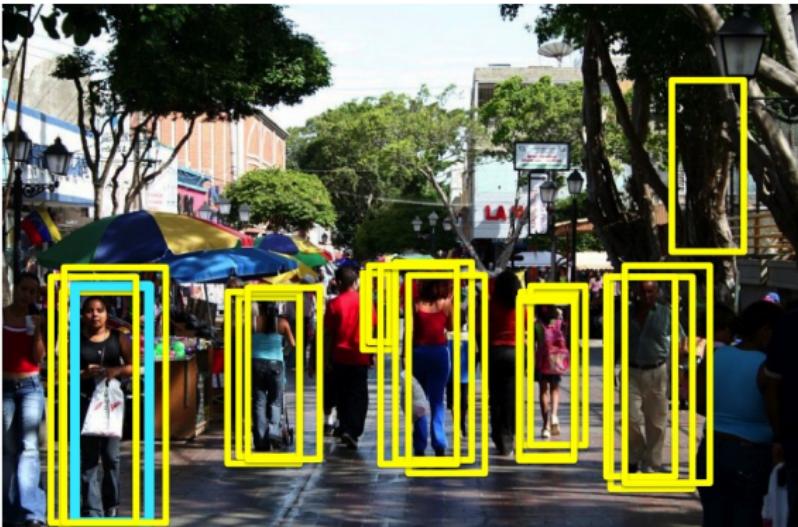
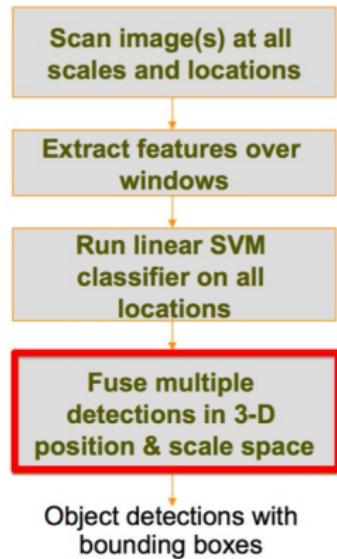


- Run detector on all scales (image sizes)
- Find scores (and thus boxes) higher than threshold
- You get a soup of overlapping boxes. What can you do to get rid of multiple detections of the same object?

The HOG Detector – Post-processing

- Perform Non-Maxima Supression (NMS)

Detection Phase



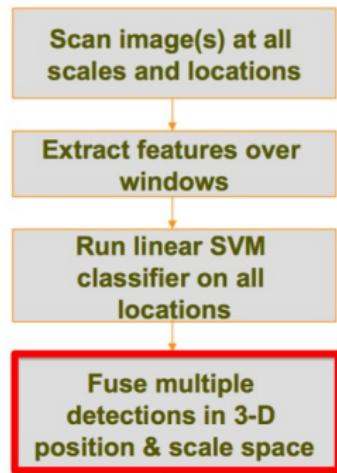
Non-maxima suppression (NMS)

- Greedy algorithm.
- At each iteration pick the highest scoring box.

The HOG Detector – Post-processing

- Perform Non-Maxima Supression (NMS)

Detection Phase



Non-maxima suppression (NMS)

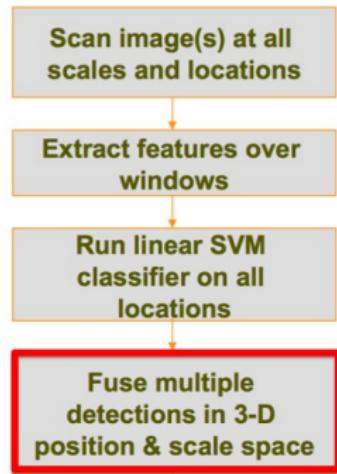
$$\text{overlap} = \frac{\text{area}(box_1 \cap box_2)}{\text{area}(box_1 \cup box_2)} > 0.5 \rightarrow \begin{array}{l} \text{remove} \\ \text{box}_2 \end{array}$$

- Remove all boxes that overlap more than XX (typically 50%) with the chosen box

The HOG Detector – Post-processing

- Perform Non-Maxima Supression (NMS)

Detection Phase



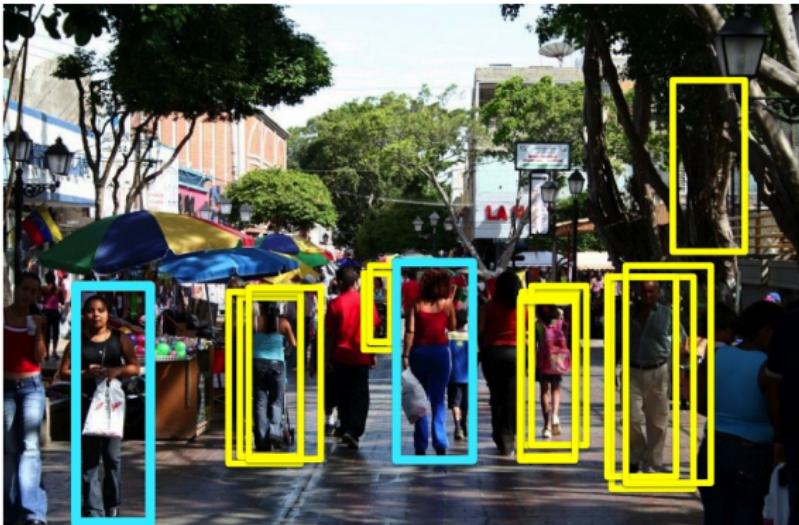
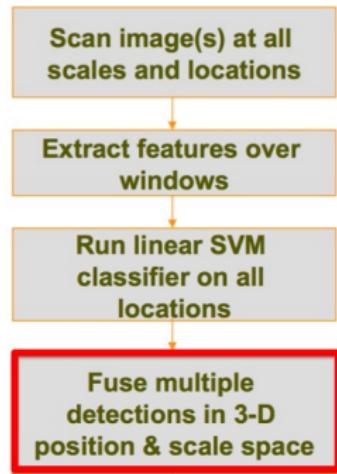
Non-maxima suppression (NMS)

- Greedy algorithm.
- At each iteration pick the highest scoring box.
- Remove all boxes that overlap more than XX (typically 50%) with the chosen box

The HOG Detector – Post-processing

- Perform Non-Maxima Supression (NMS)

Detection Phase



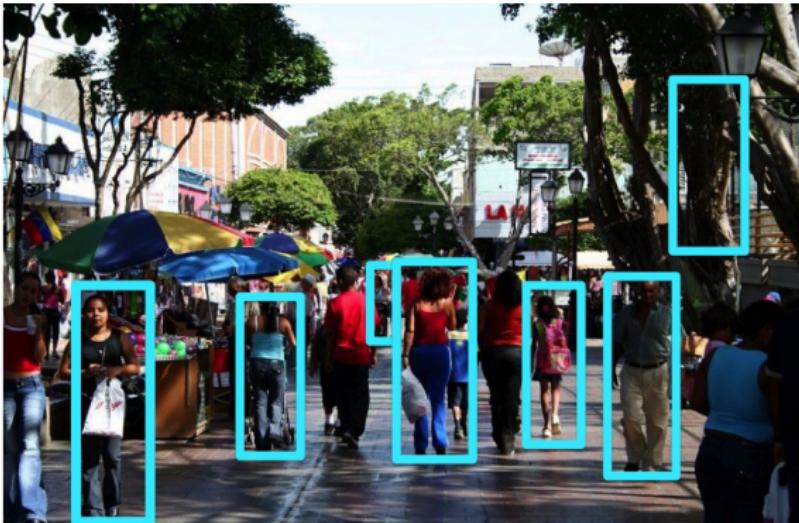
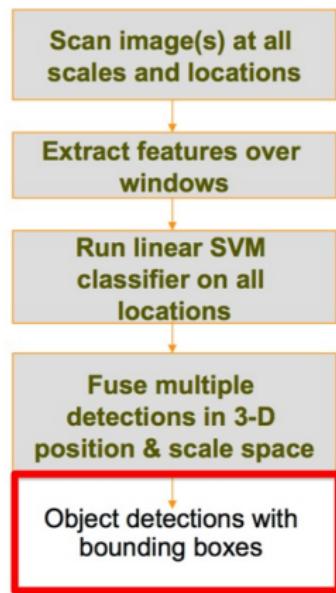
Non-maxima suppression (NMS)

- Greedy algorithm.
- At each iteration pick the highest scoring box.
- Remove all boxes that overlap more than XX (typically 50%) with the chosen box

The HOG Detector – Post-processing

- Done!

Detection Phase



Voila!

(Any idea how you would get rid of that tree detection or the upper right?)

Results

- Some results



How Should We Evaluate Object Detection Approaches?

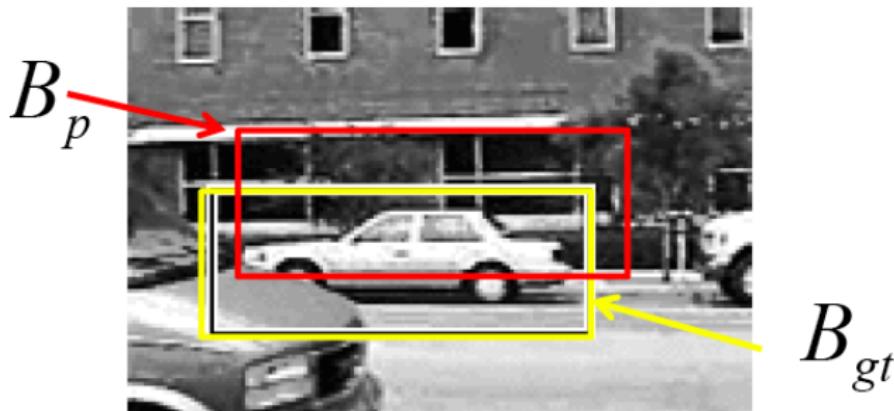
- How can we tell if our approach is doing well?
- What should be our evaluation?

What's a Correct Detection

Evaluation criteria:

- Detection is correct if the intersection of the bounding boxes, divided by their union, is > 50%.

$$a_0 = \frac{\text{area}(B_p \cap B_{gt})}{\text{area}(B_p \cup B_{gt})}$$



[Source: K. Grauman, slide credit: R. Urtasun]

Multiple Detections are Considered Wrong

- Below both detections have more than 50% overlap with ground-truth annotation. But only **one** will count as correct, the other(s) will count as **false positive** (wrong).



Precision and Recall

- We sort all the predicted boxes (for all images) according to scores, in descending order
- Then for each k we compute precision and recall obtained when using top k boxes in the list

Precision and Recall

- We sort all the predicted boxes (for all images) according to scores, in descending order
- Then for each k we compute precision and recall obtained when using top k boxes in the list
- Recall:

$$\text{recall} = \frac{\#\text{correct boxes}}{\#\text{ground-truth boxes}}$$

- Precision:

$$\text{precision} = \frac{\#\text{correct boxes}}{\#\text{all predicted boxes}}$$

- What's the min/max value of recall/precision?

Precision and Recall

- We sort all the predicted boxes (for all images) according to scores, in descending order
- Then for each k we compute precision and recall obtained when using top k boxes in the list
- Recall:

$$\text{recall} = \frac{\#\text{correct boxes}}{\#\text{ground-truth boxes}}$$

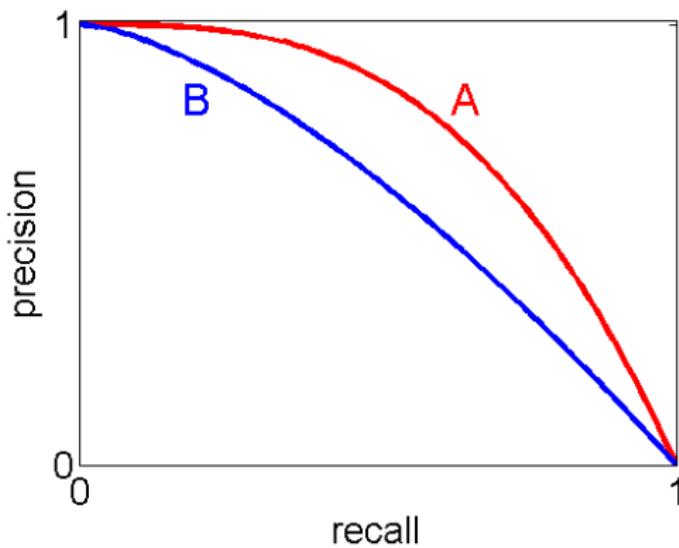
- Precision:

$$\text{precision} = \frac{\#\text{correct boxes}}{\#\text{all predicted boxes}}$$

- What's the min/max value of recall/precision?

Precision and Recall Curve

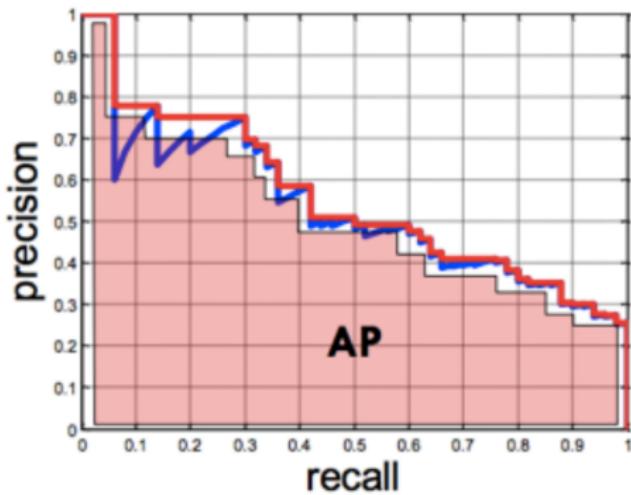
- Then you can plot a precision-recall curve
- Which curve in the plot below is better, A or B?



[Pic: http://pmtk3.googlecode.com/svn-history/r785/trunk/docs/demos/Decision_theory/PRhand_01.png]

Average Precision

- Average Precision (AP): Compute the area under the precision-recall curve
- What's the best AP one can get? What's the worst?
- AP is the standard measure for evaluating object detection performance
- Sometimes you may encounter notation mAP. This is mean Average Precision, and it's just an average of APs across different classes.



[Pic from: R. Girshik]

Performance of the HOG Detector (back in 2005)

- Interesting: Look at the curve for PCA-SIFT (improved SIFT). Way down there...
- Note: Be careful. HOG performs better for this task, a task and data set it was specifically tuned for, and the performance gap is within the variation caused by that tuning.
- Other datasets with different properties, for example where rotation in-variance is crucial will yield different results.

