

# Detection & recognition



한양대학교  
HANYANG UNIVERSITY

Tae Hyun Kim



# Roadmap

## Classification



CAT

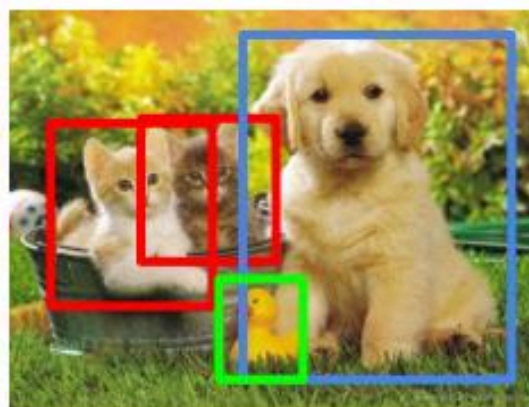
## Classification + Localization

→ 위치 찾음



CAT

## Object Detection



CAT, DOG, DUCK

## Instance Segmentation



CAT, DOG, DUCK

Single object

Multiple objects



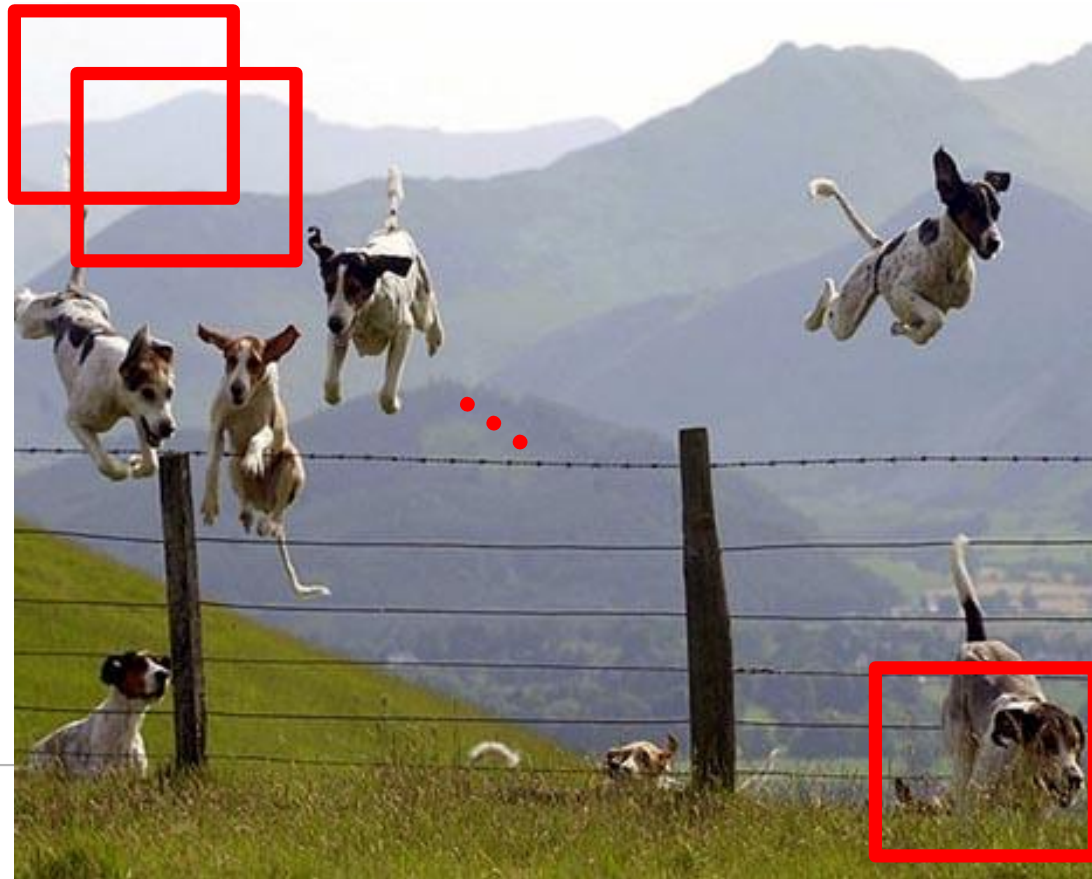
# State-of-the-art Detection Method

- YOLO v2 (CVPR'18)

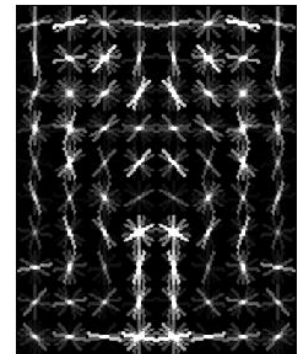


# Object Category Detection

- Focus on object search: “Where is it?”
- Build templates that quickly differentiate object patch from background patch



Dog Model



**Object** or  
**Non-Object?**





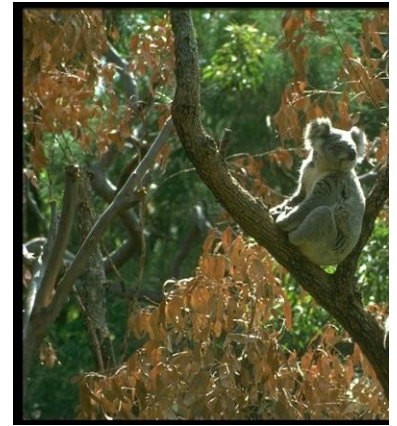
# Challenges in modeling the object class



Illumination



Object pose



Clutter

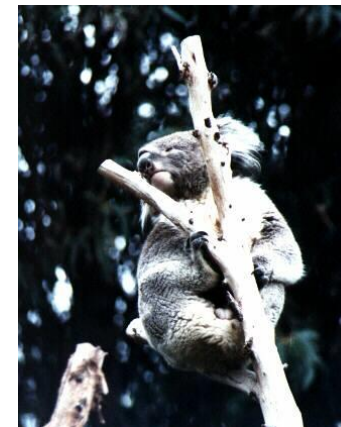
$m(\psi_x \text{ variation}) \propto \frac{3}{2}$



Occlusions



Intra-class  
appearance



Viewpoint



# Challenges in modeling the non-object class

True Detections



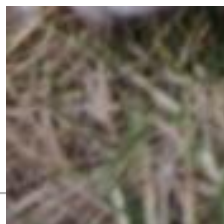
Bad Localization



Confused with Similar Object



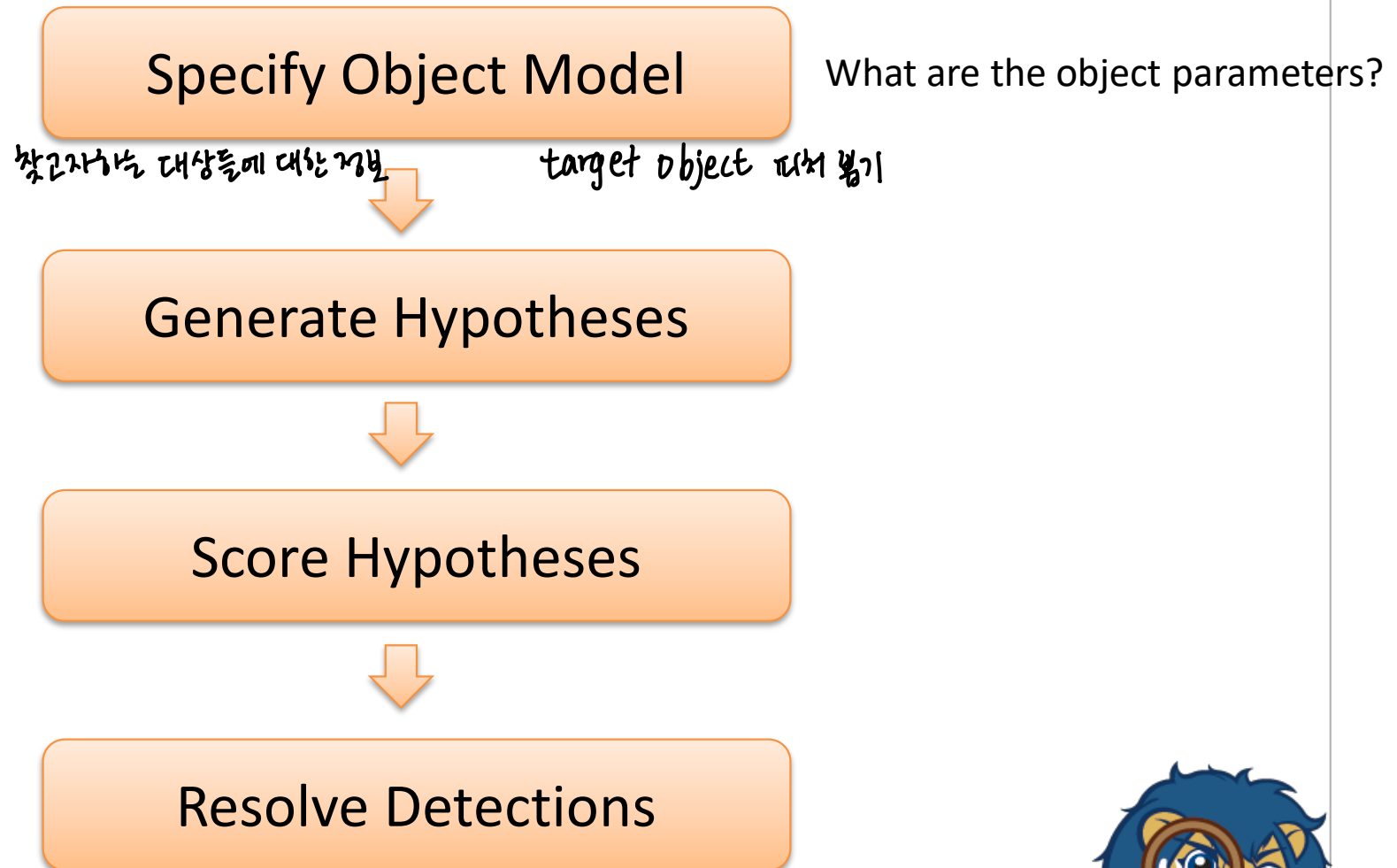
Misc. Background



Confused with Dissimilar Objects



# General Process of Object Detection





# Specifying an object model

## 1. Statistical Template in Bounding Box

- Object is somewhere  $(x,y,w,h)$  in image
- Features defined w.r.t. bounding box coordinates

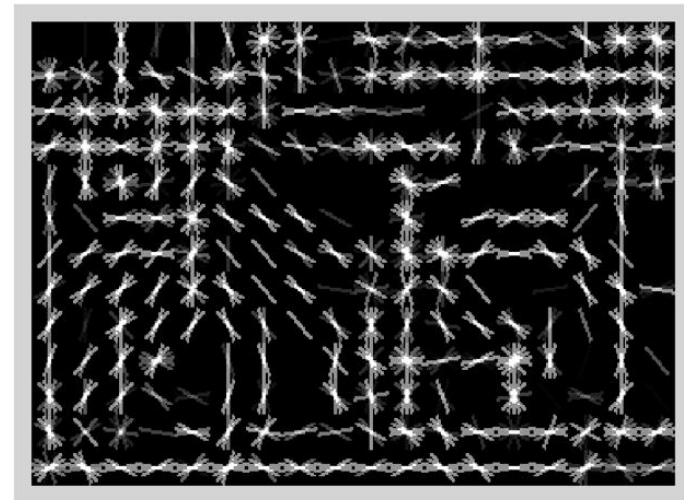
찾고자 하는 객체의 피쳐화

'좋은 피쳐' - 같은 카테고리 안에서 공통적인 피쳐

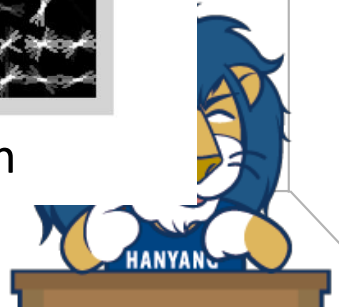


Image

gradient 정보



Template Visualization



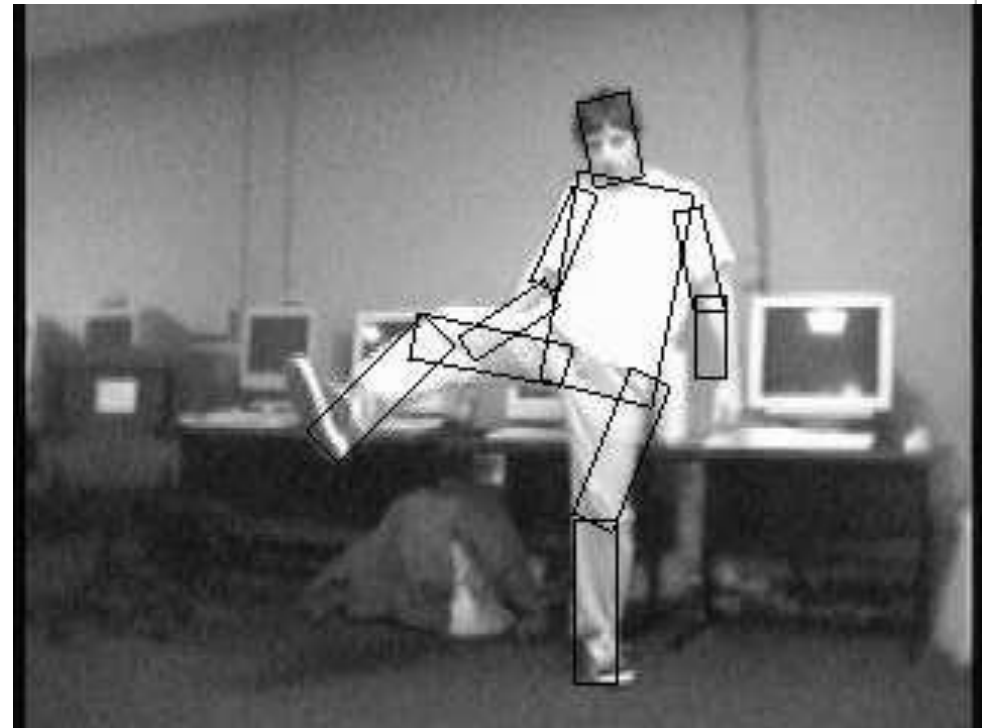
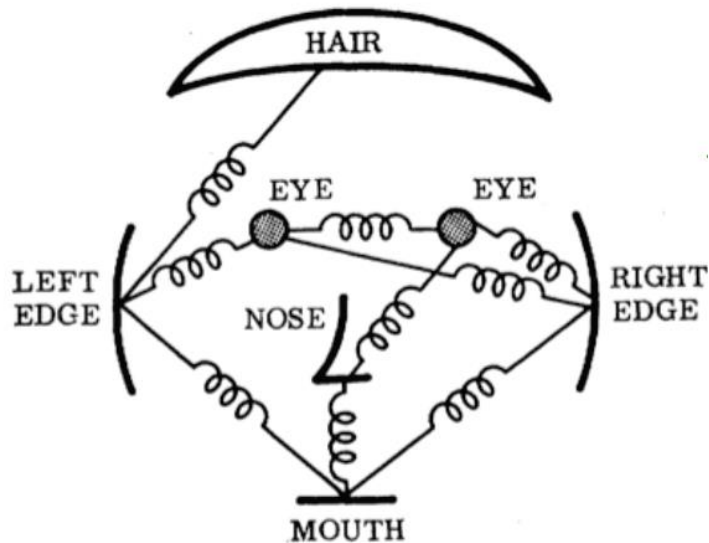


# Specifying an object model

## 2. Articulated parts model

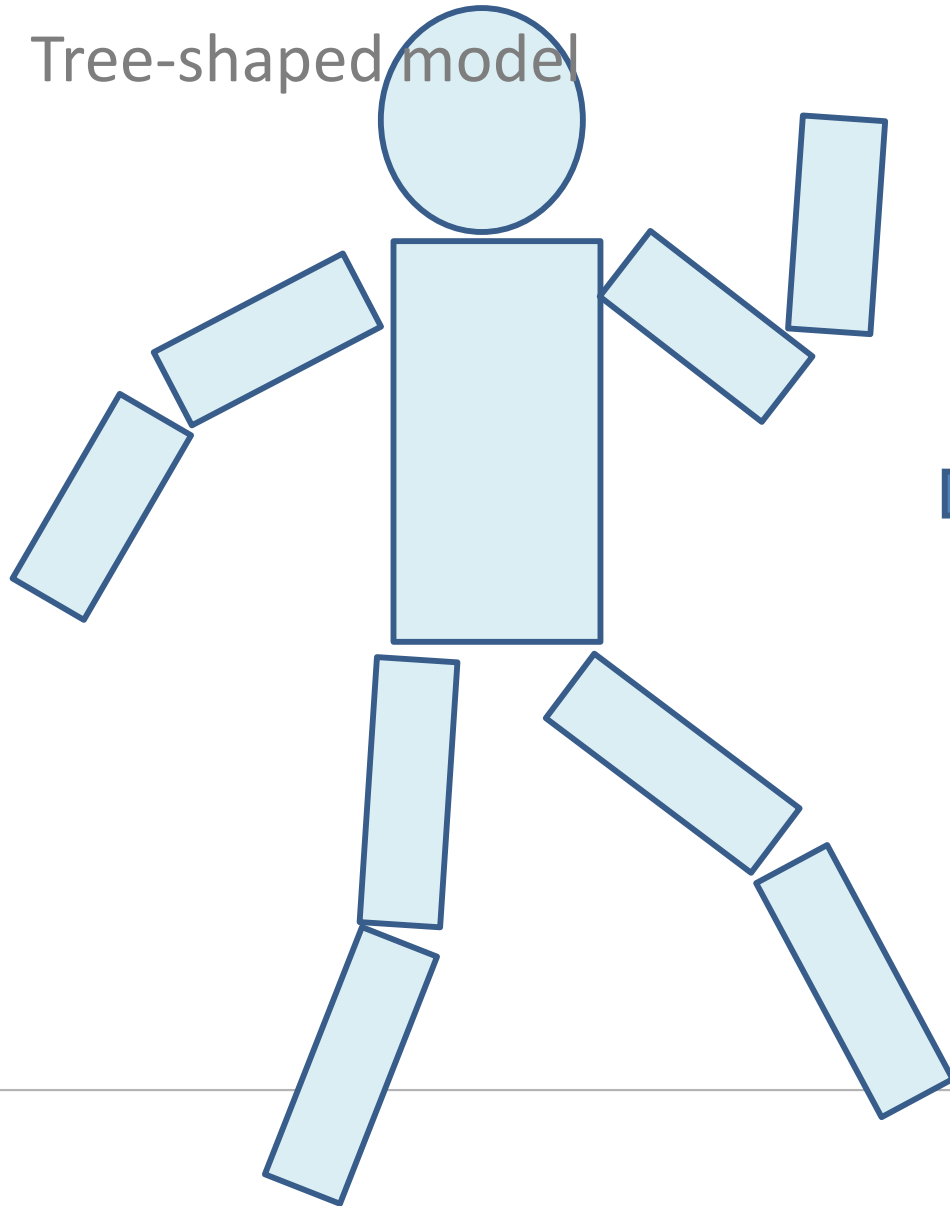
- Object is configuration of parts
- Each part is detectable

사진 constraint를 가지고  
object define

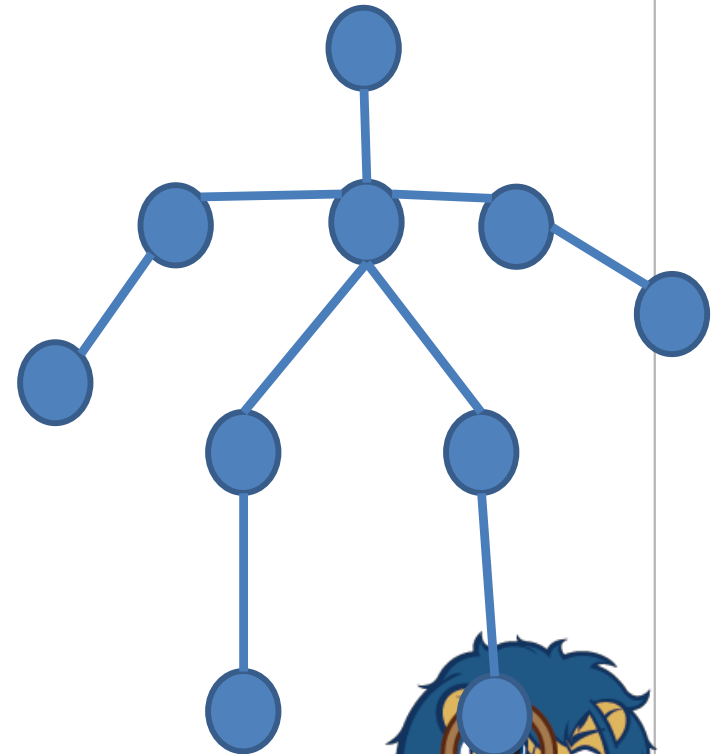


# How to model spatial relations?

- Tree-shaped model

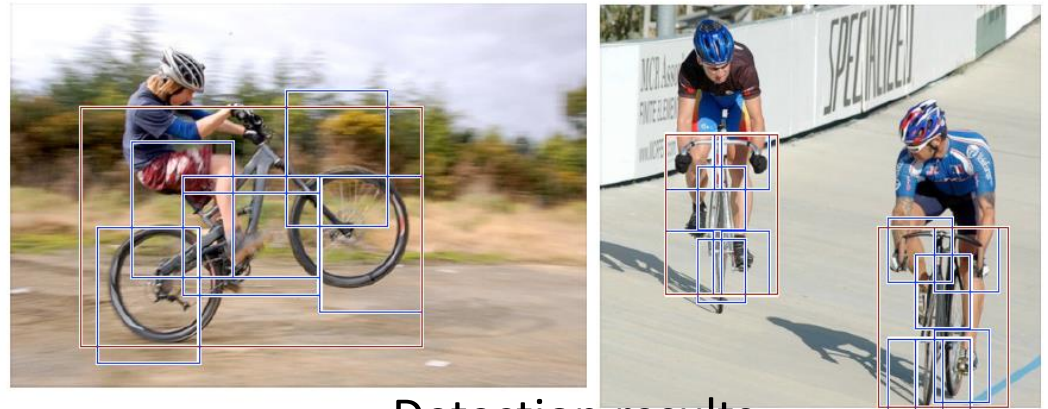


사전 정보로 더 좋은 피쳐 만들기



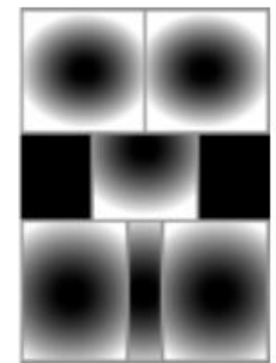
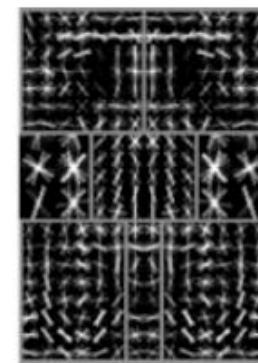
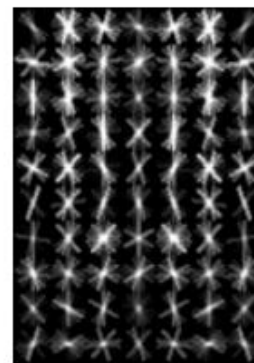
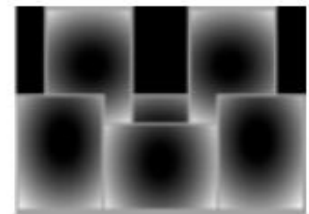
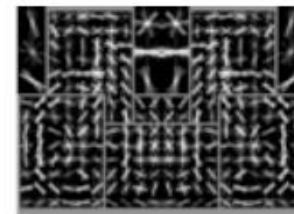
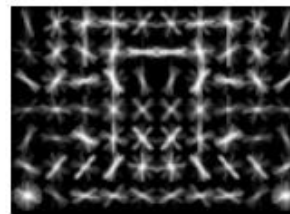
# Specifying an object model

## 3. Hybrid template/ parts model



Detection results

## Template Visualization

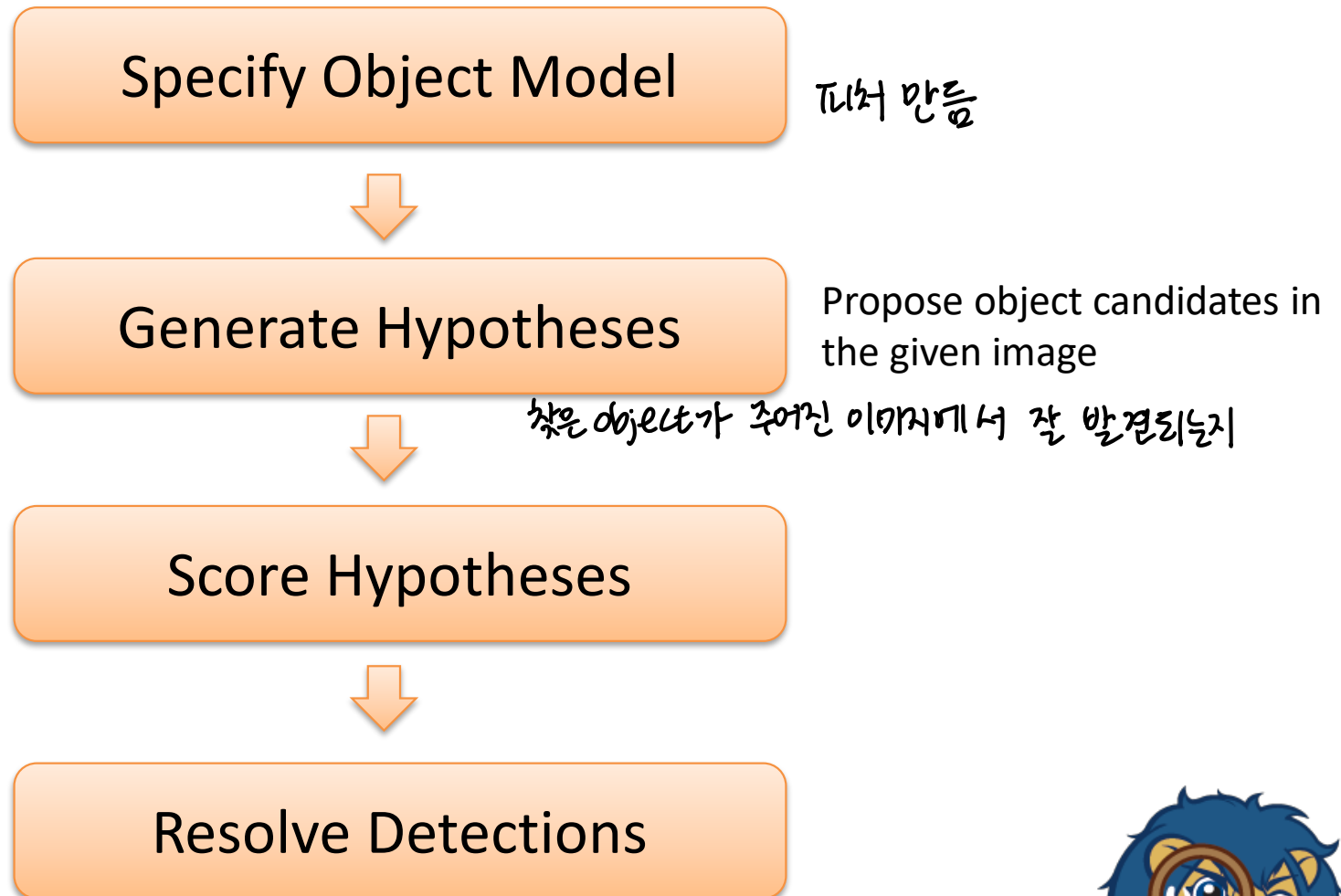


root filters  
coarse resolution

part filters  
finer resolution

deformation  
models

# General Process of Object Detection





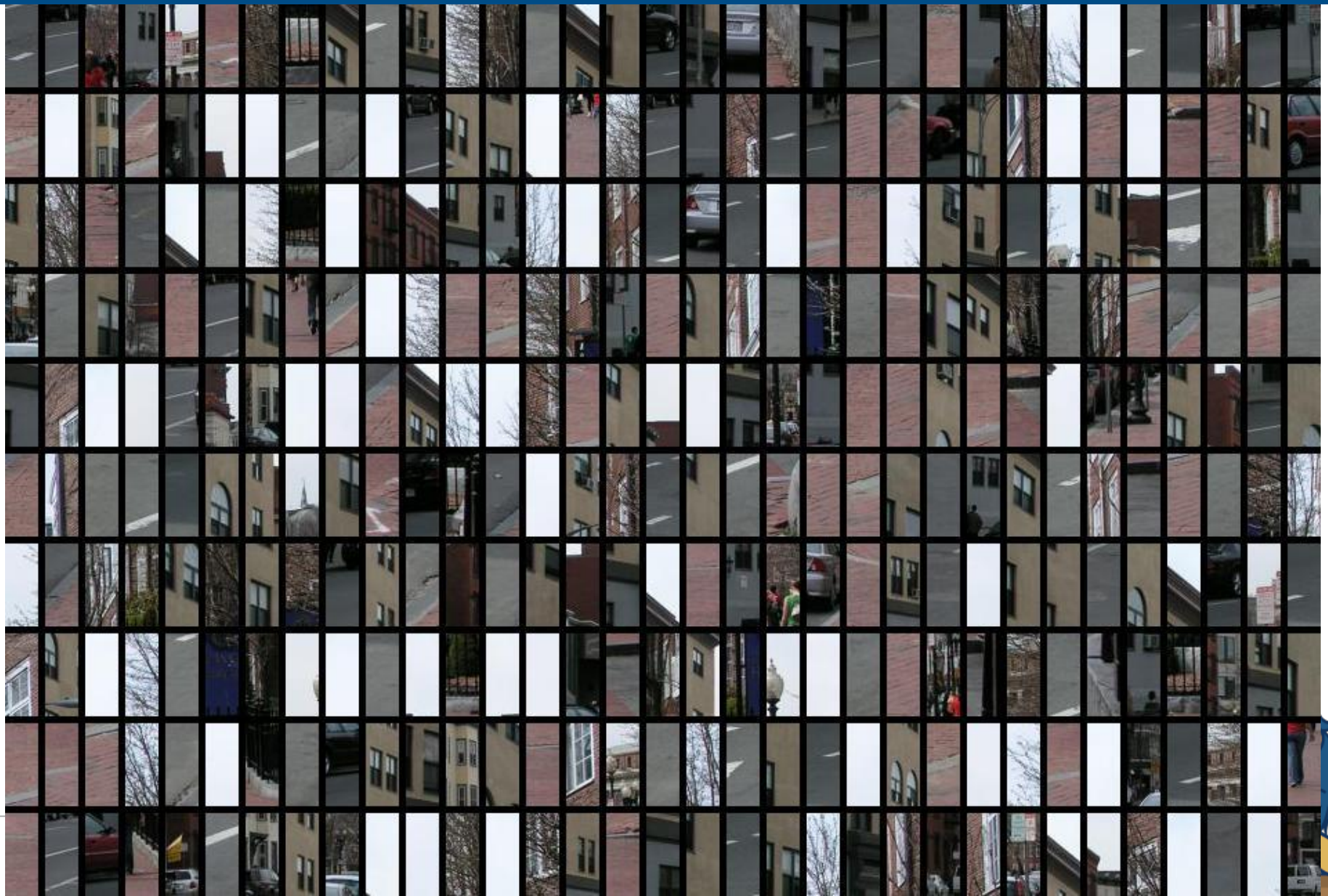
# Generating hypotheses

1. Sliding window
  - Test patch at each location and scale



# Generated hypotheses

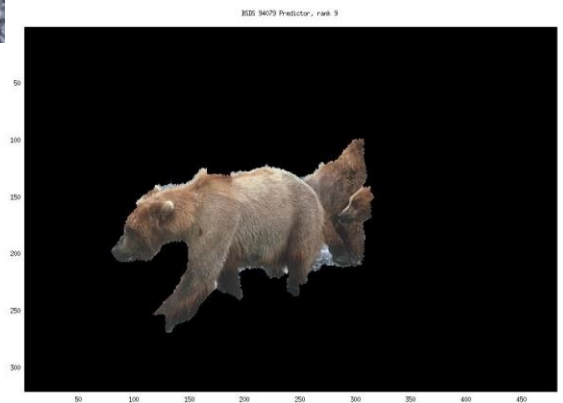
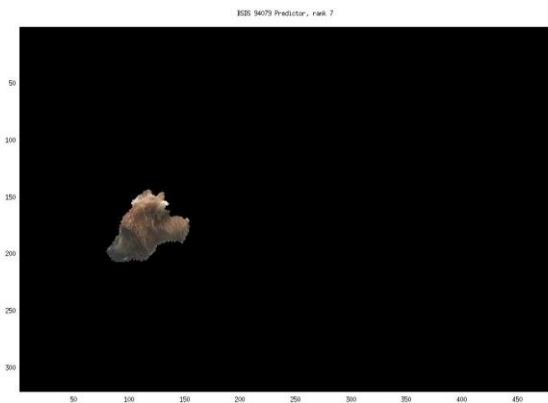
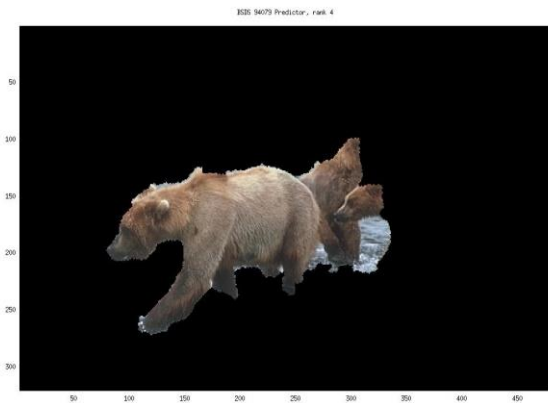
candidate → object 유무 검사



# Generating hypotheses

## Region-based proposal

비슷한 컬러끼리 뮌쳐닝





# General Process of Object Detection

Specify Object Model



Generate Hypotheses

후보 생성



Score Hypotheses

Compare hypotheses and the  
trained models (templates).  
Traditional classifiers/Neural nets

피처의 유사도 측정  
Similarity



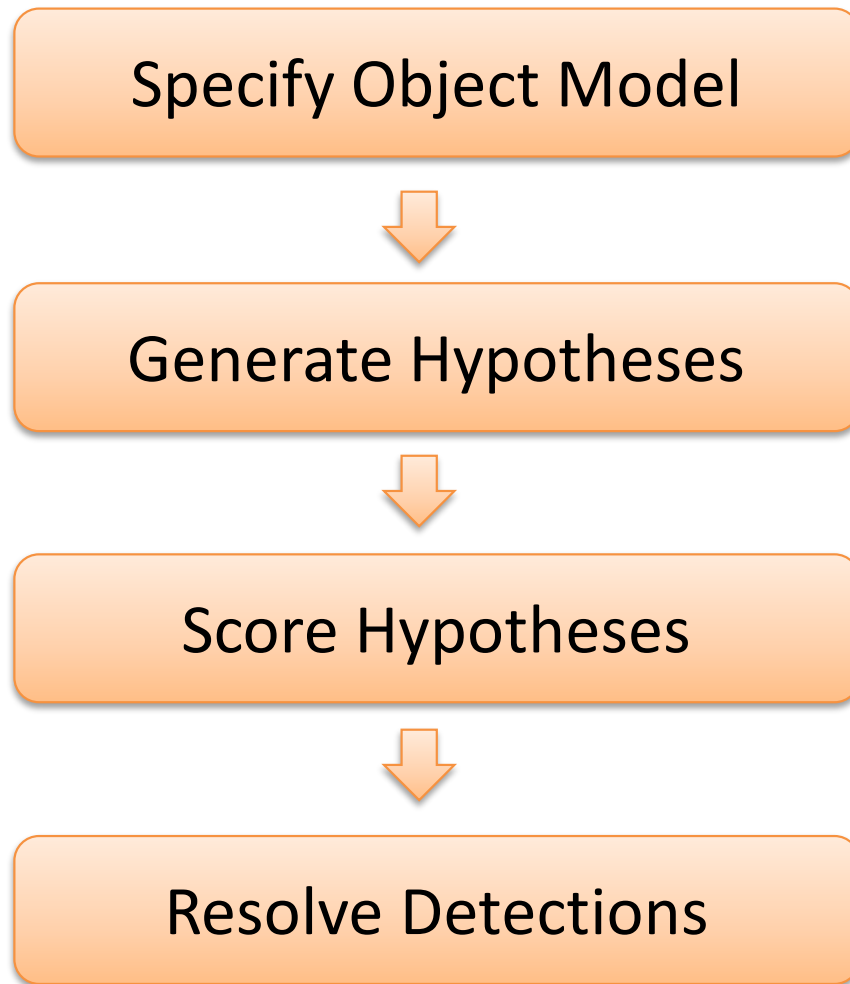
ex) 유클리디안 distance

Resolve Detections





# General Process of Object Detection



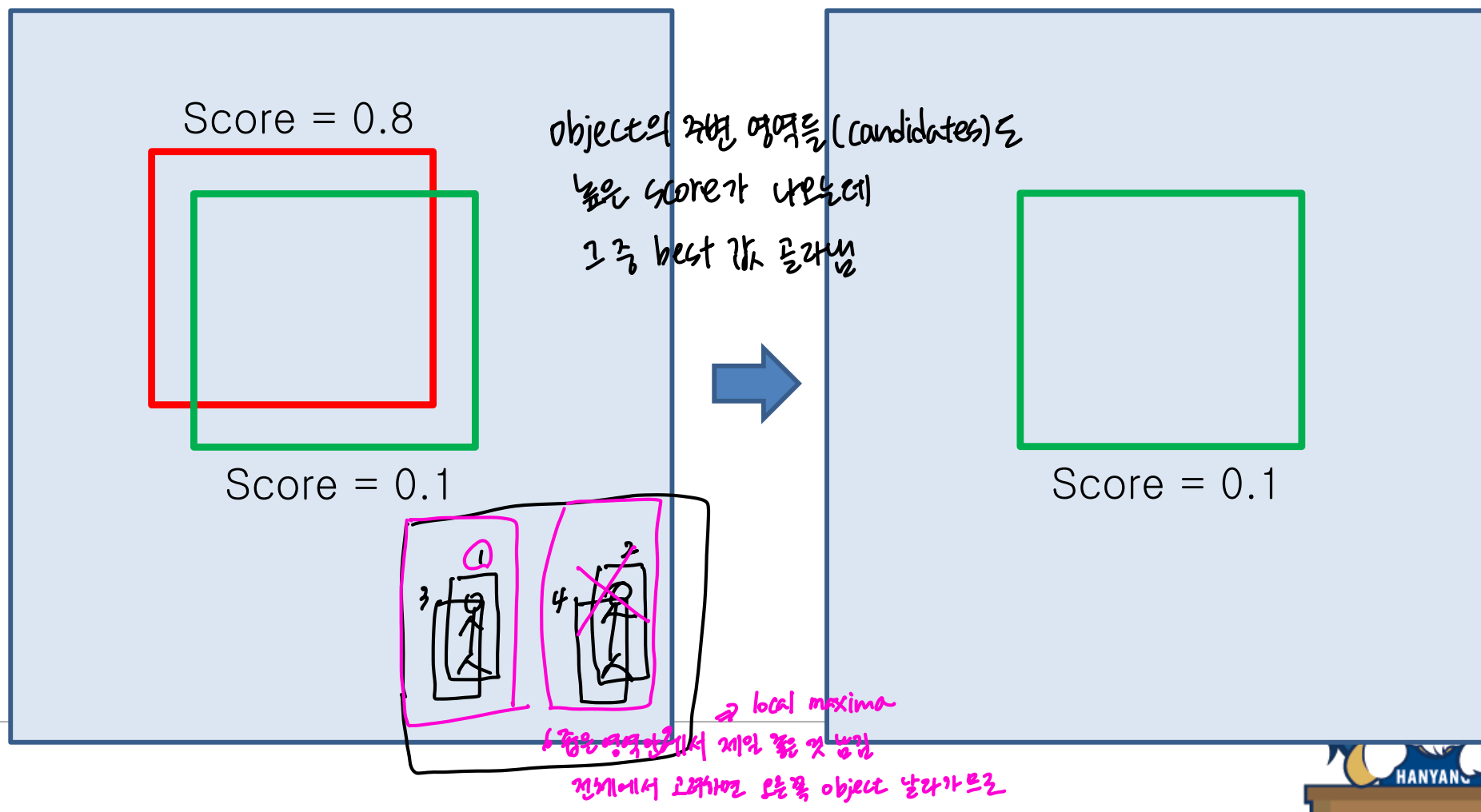
피처가 optimal 하지 않음

피처의 퀄리티 차이



# Resolving detection scores

- Non-max suppression (NMS) 꼭 필요한 단계
  - Find local maxima

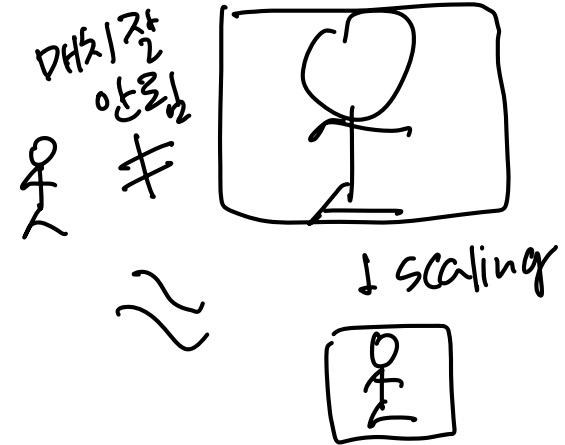


# Object category detection in computer vision

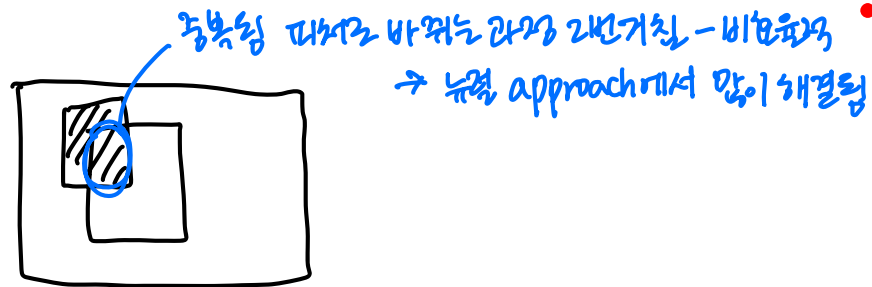
- Goal: detect all pedestrians, cars, monkeys, etc in image



# Sliding window: a simple alignment solution



모든 대치들에 대해  
사람이 들어있는지 검사



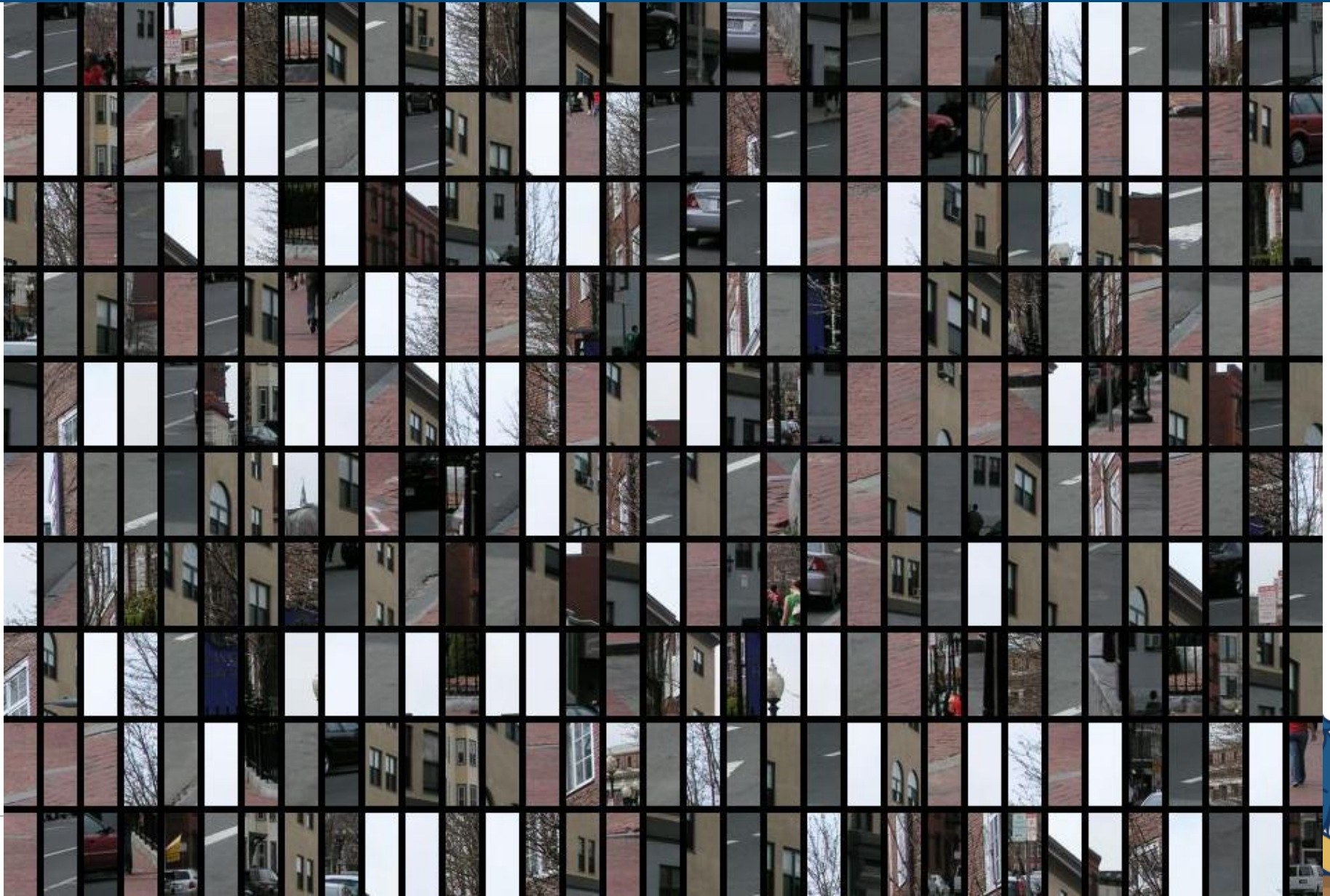
동영상 피쳐를 바꿔는 과정 2번 거친 - 비효율적  
→ 뉴럴 approach에서 많이 해결됨





# Each window is separately classified

각 candidate에서 피쳐 뽑아서 유사도 측정



# Design challenges

- How to efficiently search for likely objects
  - Even simple models require searching *hundreds of thousands of positions and scales* 스케일 고려
- Feature design and scoring
  - How should *appearance* be modeled? What features correspond to the object? 외형 무시하고 같은 카테고리라는 것을 만족시키는 피쳐
- How to deal with different viewpoints?
  - Often train different models for a few *different viewpoints*
- Implementation details
  - Window size
  - Aspect ratio
  - Translation/scale step size
  - Non-maxima suppression



## Example: Dalal-Triggs detector

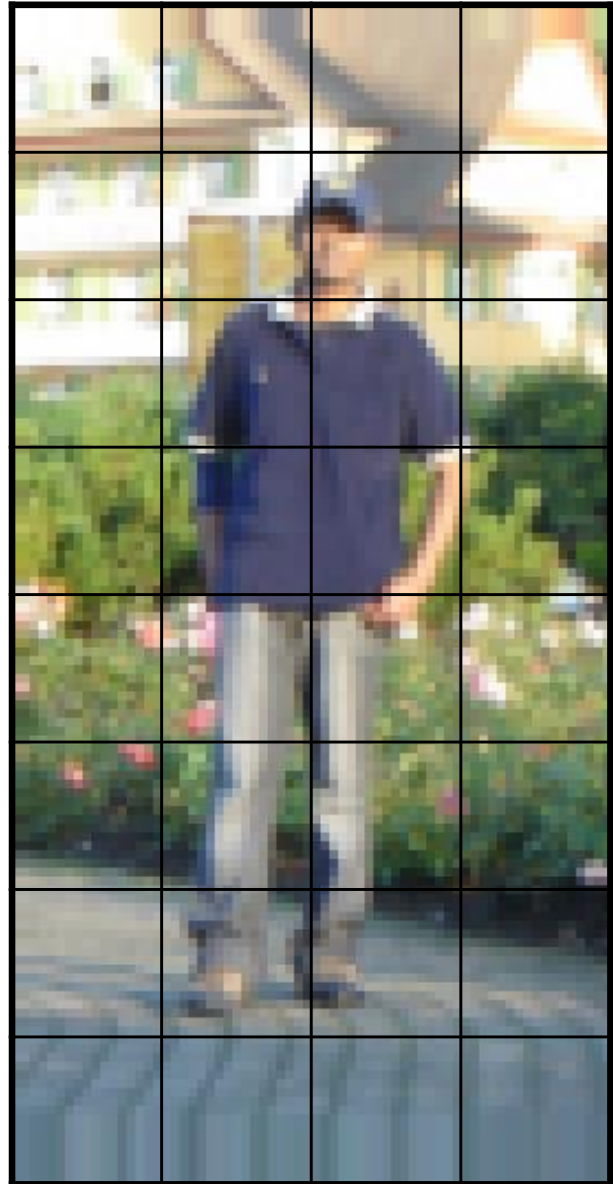
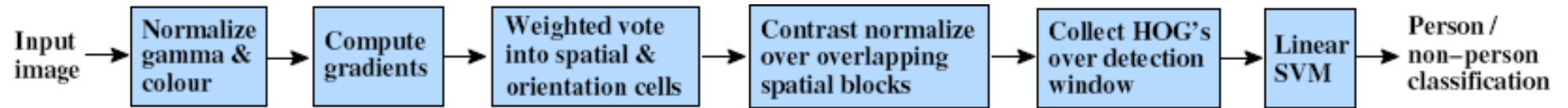


1. Extract fixed-sized (64x128 pixel) window *at each position and scale*  
candidates 만들
2. Compute **HOG** (histogram of gradient) features within each window  
각 candidates에서 피쳐 추출
3. Score the window with a linear **SVM classifier**  
내가 가지고 있는 피쳐와 유사도 측정
4. Perform **non-maxima suppression** to remove overlapping detections with lower scores

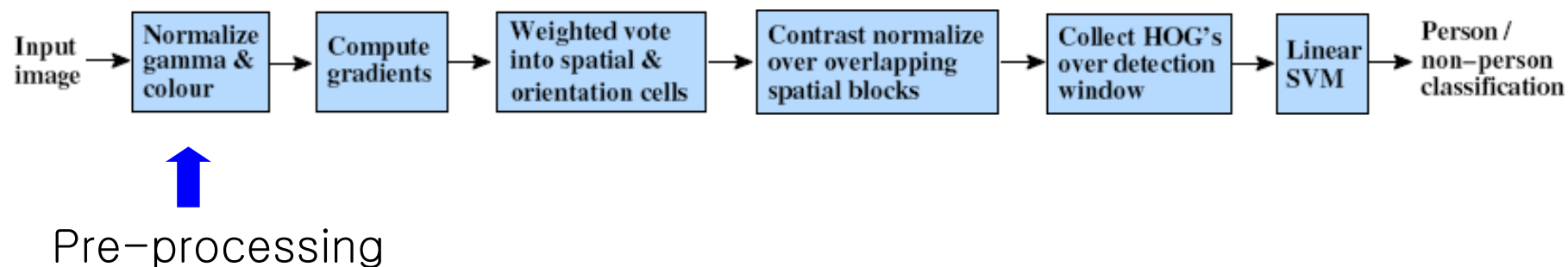




## < feature extraction overall flow >



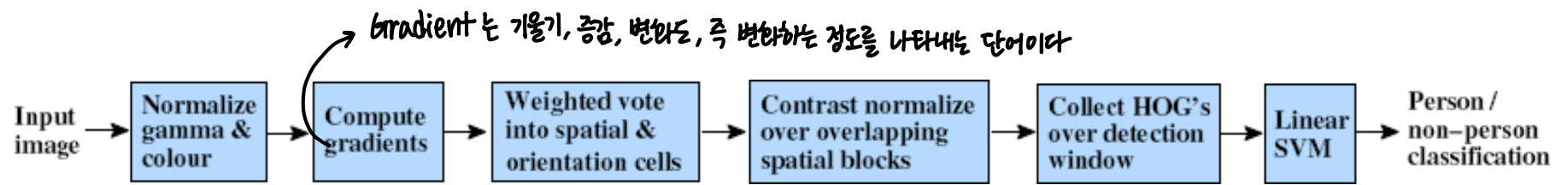




- Tested with
  - RGB
  - LAB
  - Grayscale

} Slightly better performance vs. grayscale
- Gamma Normalization and Compression
  - Square root
  - Log

} Very slightly better performance vs. no adjustment



Outperforms

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

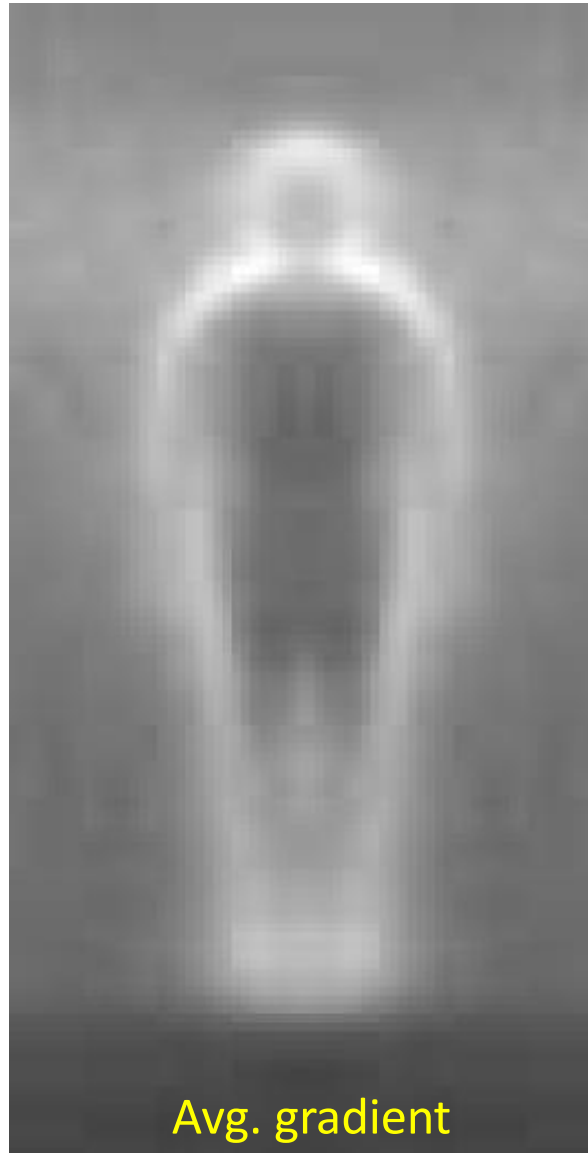
centered

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

uncentered

$$\begin{bmatrix} 1 & -8 & 0 & 8 & -1 \end{bmatrix}$$

cubic-corrected



gradient 정보 활용

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

diagonal

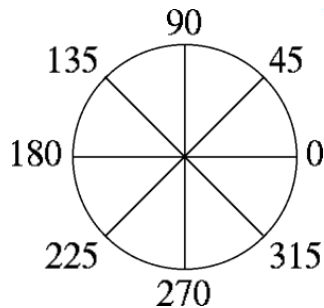
$$\begin{bmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{bmatrix}$$

Sobel

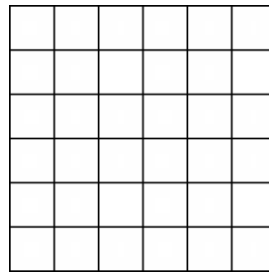


- Histogram of gradient orientations

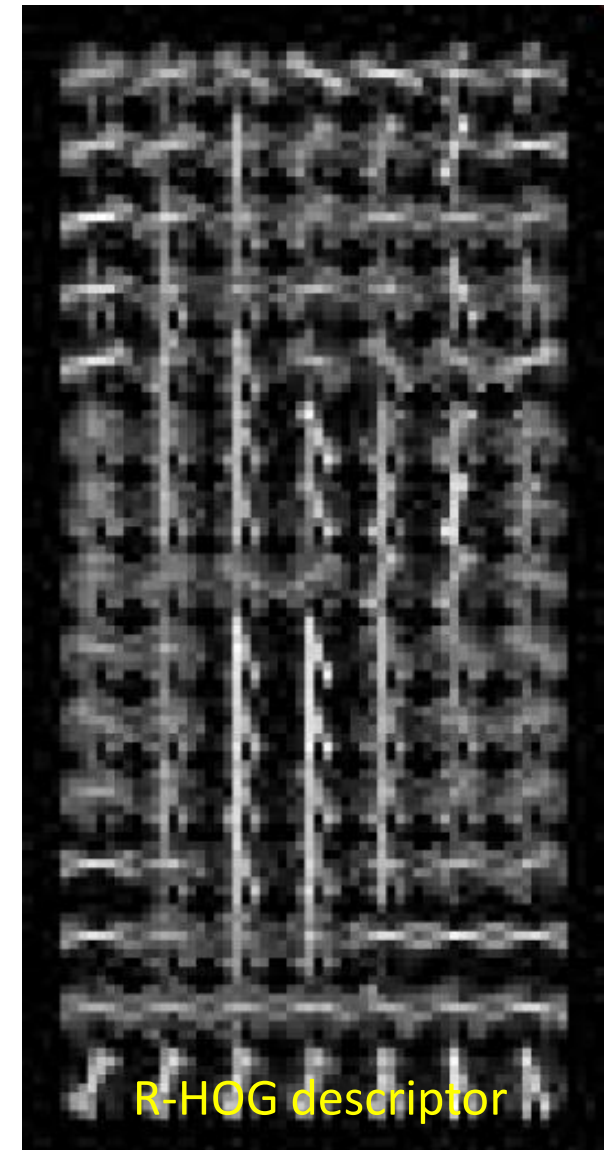
Orientation: 9 bins (for unsigned angles)

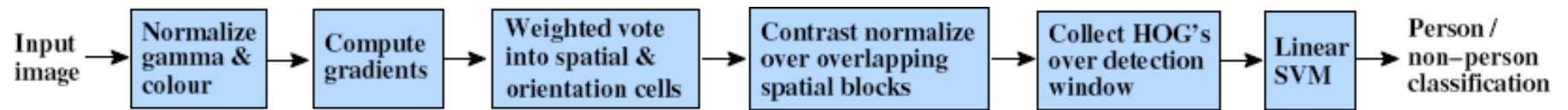


Histograms in 8x8 pixel cells



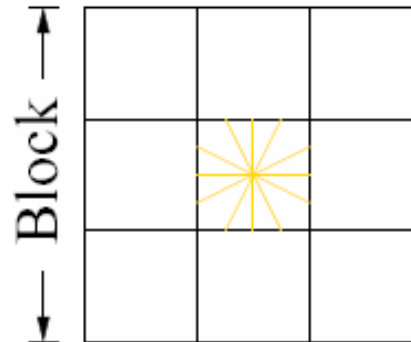
- Votes weighted by magnitude
- Bilinear interpolation between cells





## R-HOG

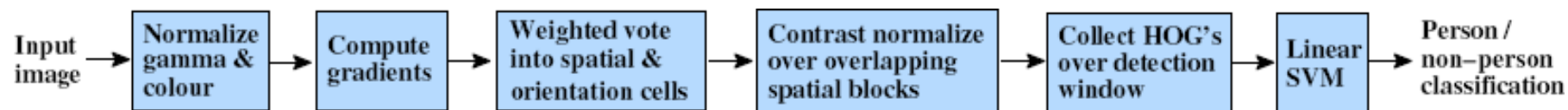
Cell



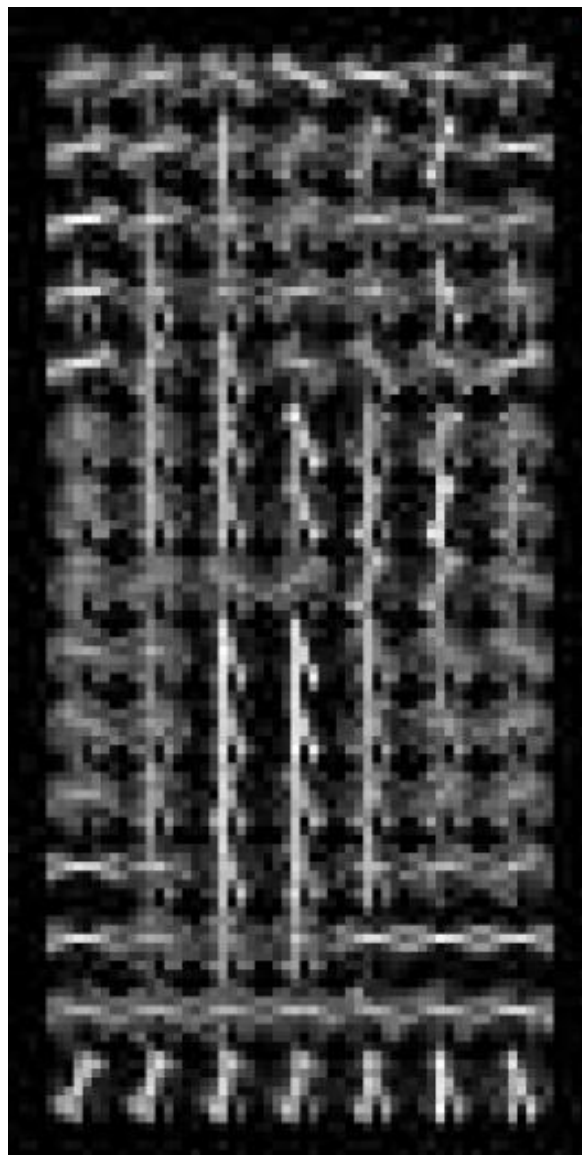
Normalize with respect to surrounding cells

$$L2 - norm : v \longrightarrow v / \sqrt{\|v\|_2^2 + \epsilon^2}$$





X=

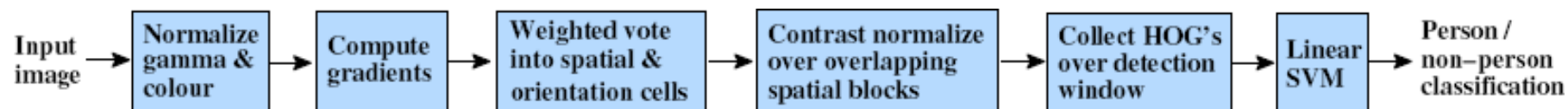


# orientations

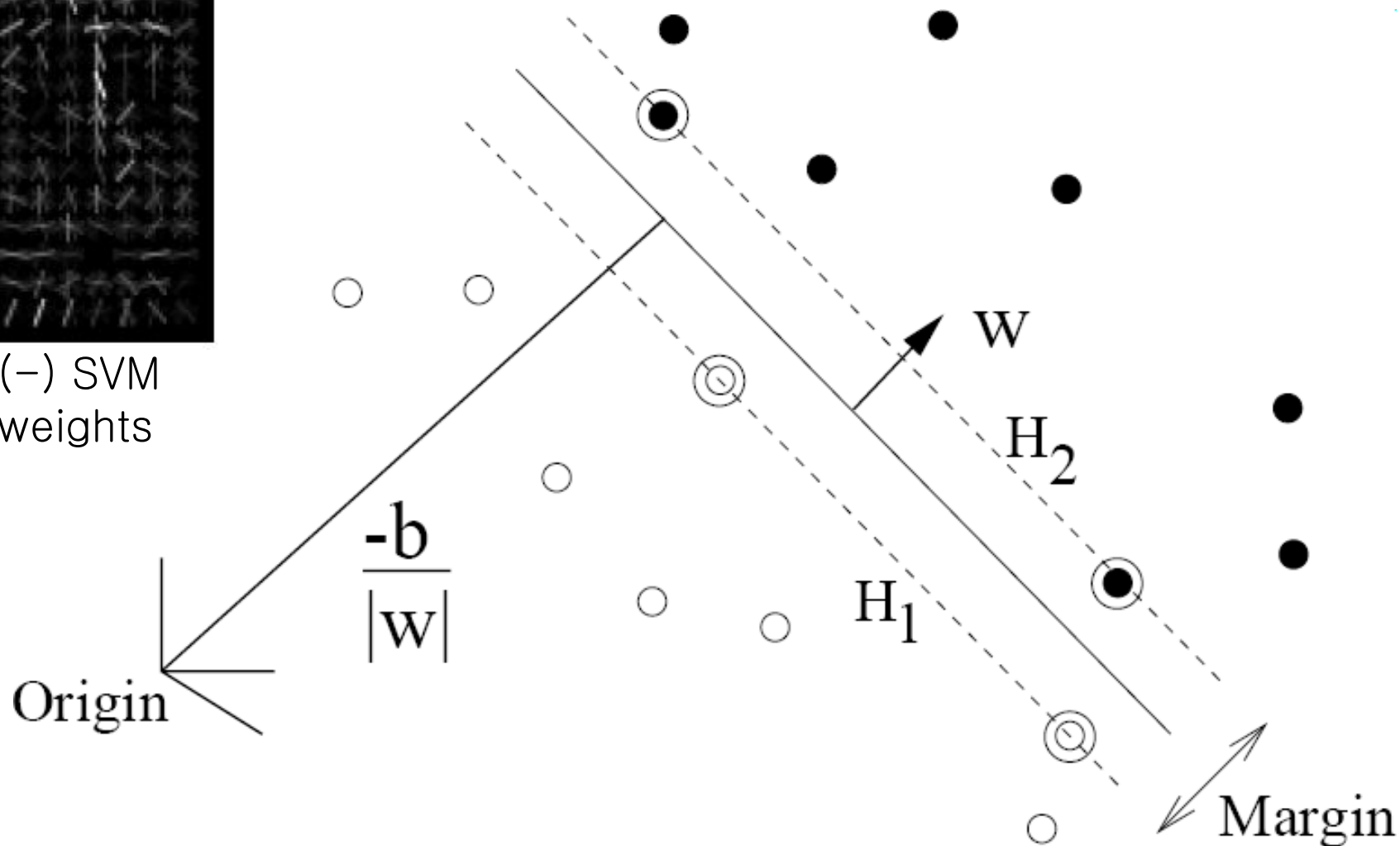
# features = 15 x 7 x 9 x 4 = 3780

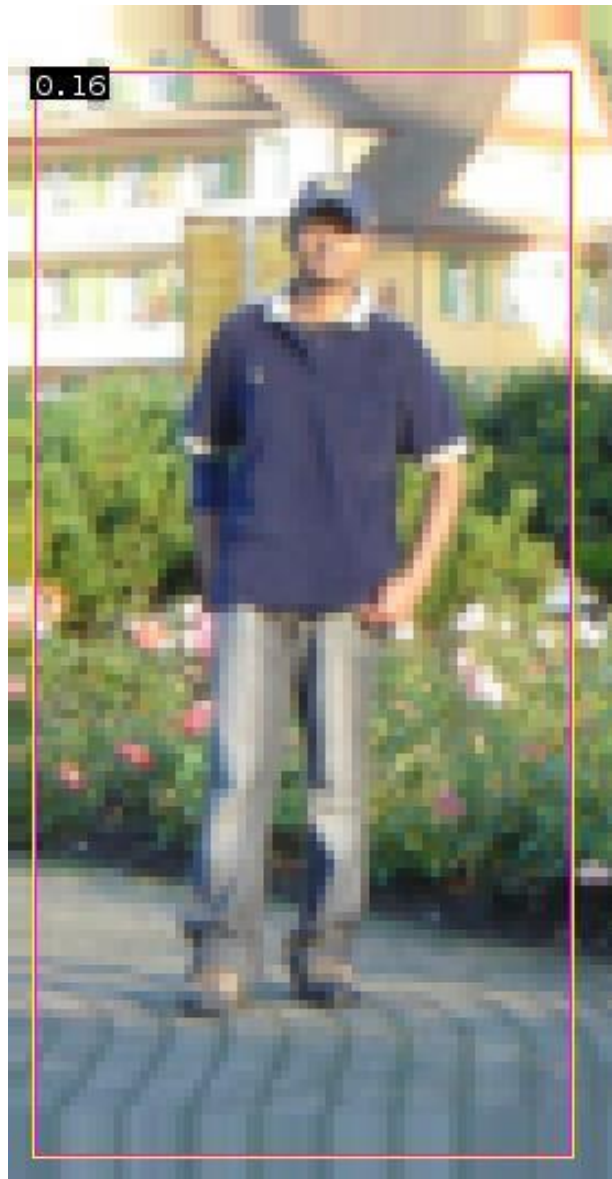
# cells

# normalizations by neighboring cells



↑  
classifier





$$0.16 = w^T x - b$$

$$\text{sign}(0.16) = 1$$

$\Rightarrow$  pedestrian

# Detection examples





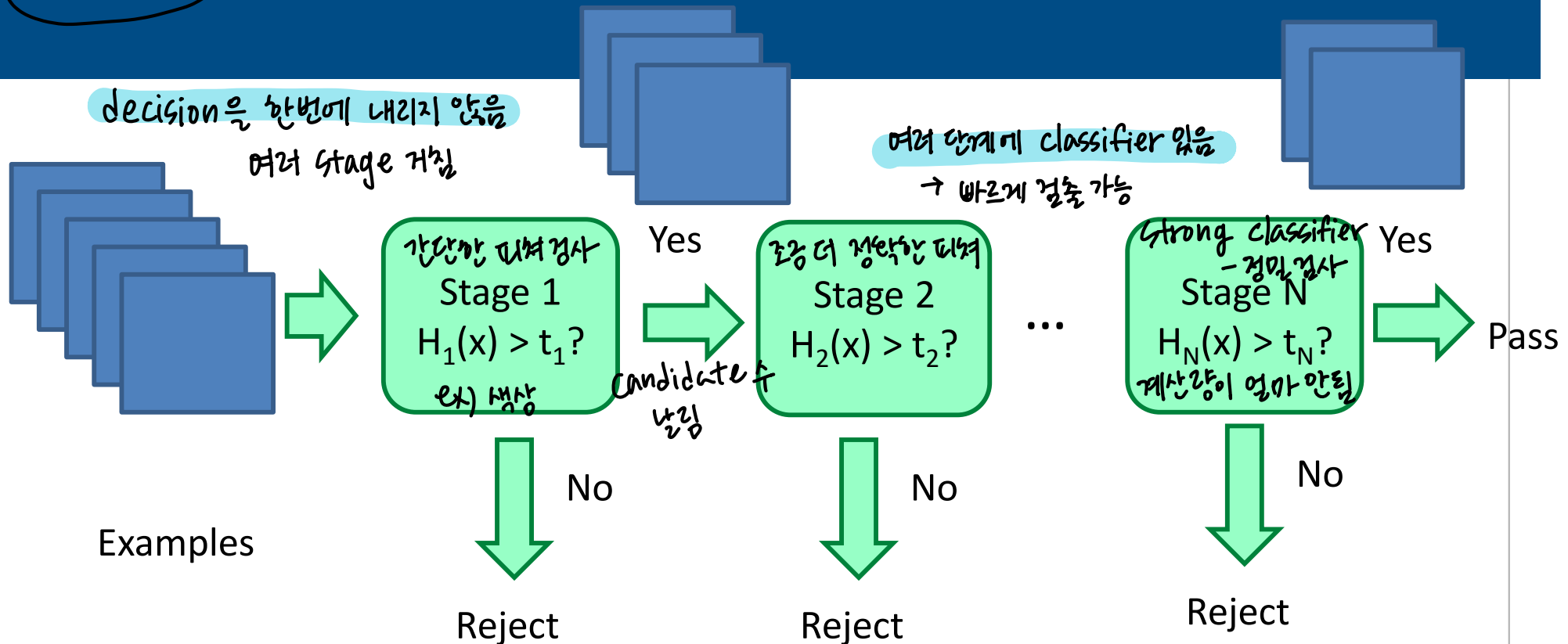
# Viola-Jones sliding window detector

*Fast detection* through two mechanisms

- Quickly *eliminate unlikely windows*
- Use *features that are fast* to compute



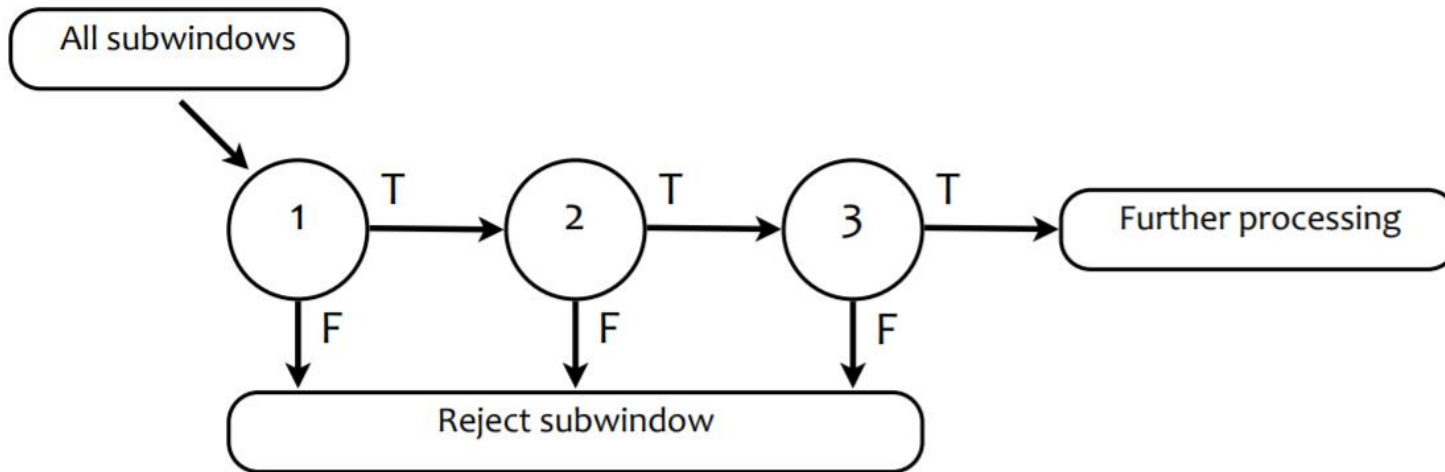
# Cascade for Fast Detection



- Choose threshold for low false negative rate
  - Avoid rejecting facial images
- Fast classifiers early stage
- Slow classifiers later, but most examples don't get there



# Cascaded Classifier



The cascaded classifier is almost *10 times faster* since only positive examples are considered in the subsequent stages

- We need to determine
  - The number of stages
  - Weak classifiers (features) in each stage.
  - The threshold of each stage
- Finding the optimal parameter is extremely difficult.



# Characteristics of AdaBoost

- Training AdaBoost is the procedure to find
  - A **set of weak classifiers** with associated features
  - The **weight** of each weak classifier
- A collection of weak classifiers
  - Each **rectangle feature** is used to construct a weak classifier
  - The strong classifier is simply a set of weak classifiers
- An iterative training procedure
  - AdaBoost performs a series of trials.
  - A new weak classifier with associated weight is selected in a greedy manner in each iteration
  - Jointly find weak classifiers and extract features

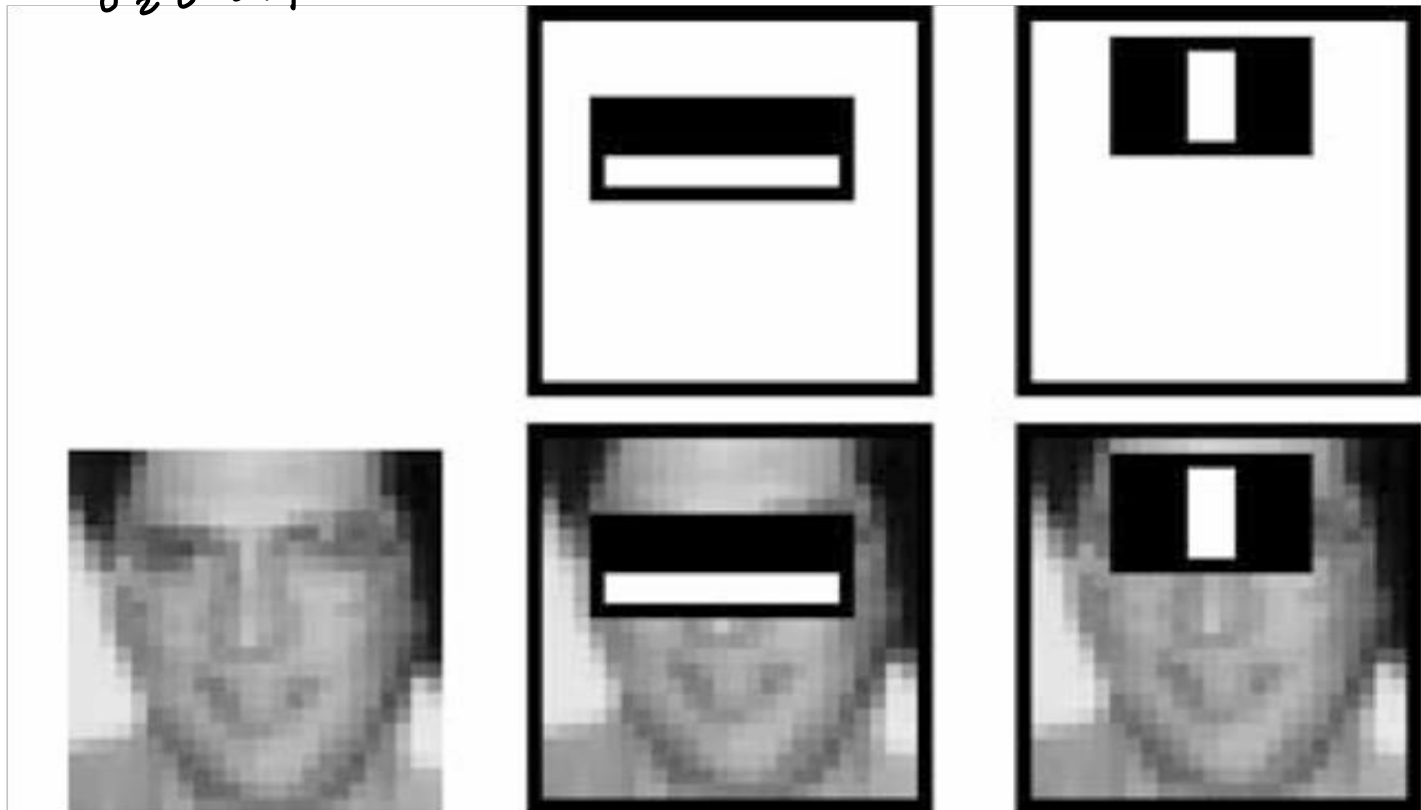




## Top 2 selected features

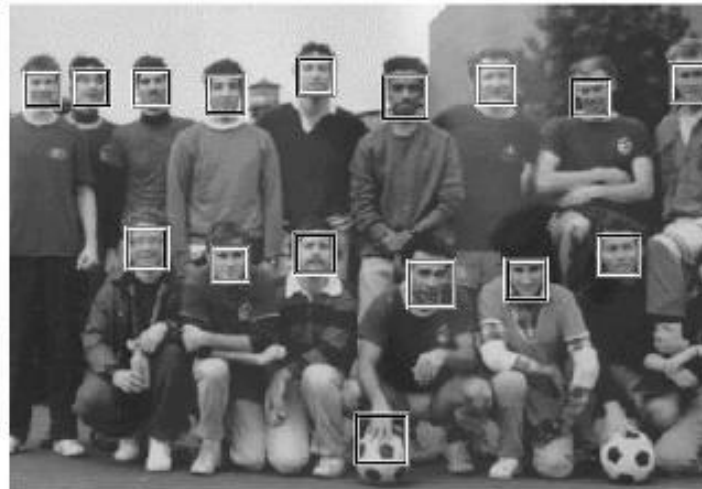
- These two-feature classifier can be adjusted to detect 100% of the faces with a false positive rate of 50%

심플한 피쳐



# Viola Jones Results

Speed = 15 FPS (in 2001)

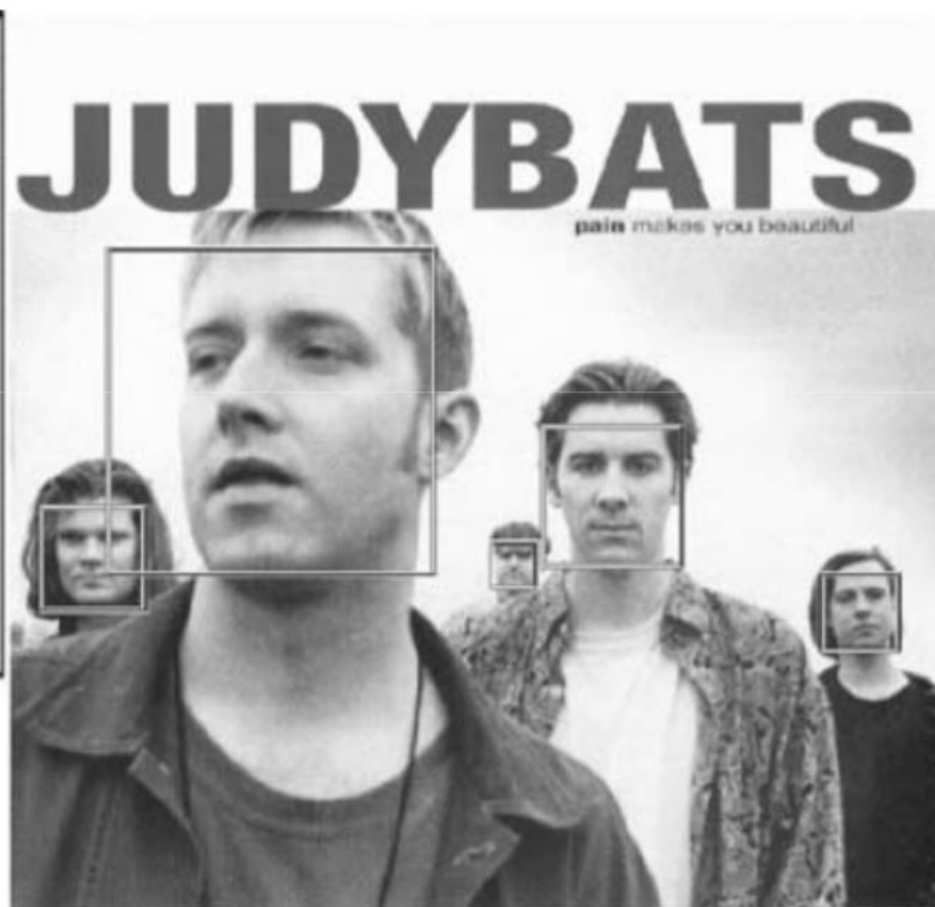


<div>False detections</div> <div>Detector</div>	10	31	50	65	78	95	167
Viola-Jones	76.1%	88.4%	91.4%	92.0%	92.1%	92.9%	93.9%
Viola-Jones (voting)	81.1%	89.7%	92.1%	93.1%	93.1%	93.2 %	93.7%
Rowley-Baluja-Kanade	83.2%	86.0%	-	-	-	89.2%	90.1%
Schneiderman-Kanade	-	-	-	94.4%	-	-	-
Roth-Yang-Ahuja	-	-	-	-	(94.8%)	-	-

MIT + CMU face dataset



# Viola Jones Results



# Thank you!

