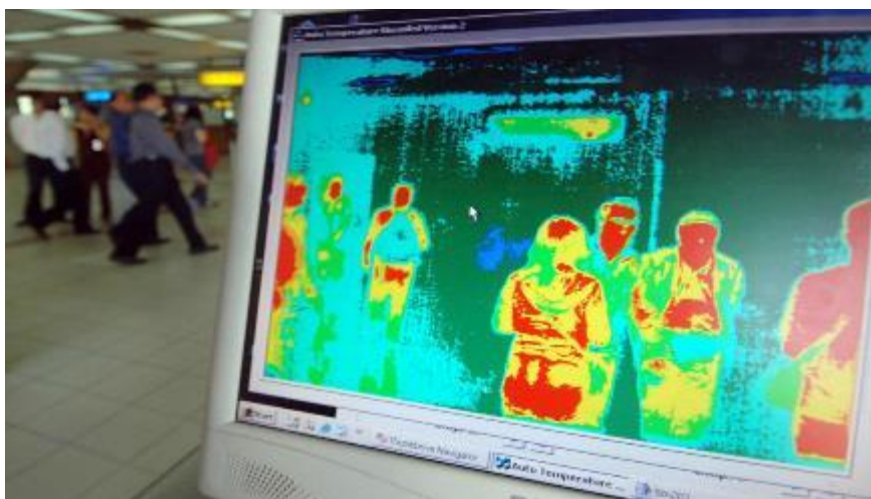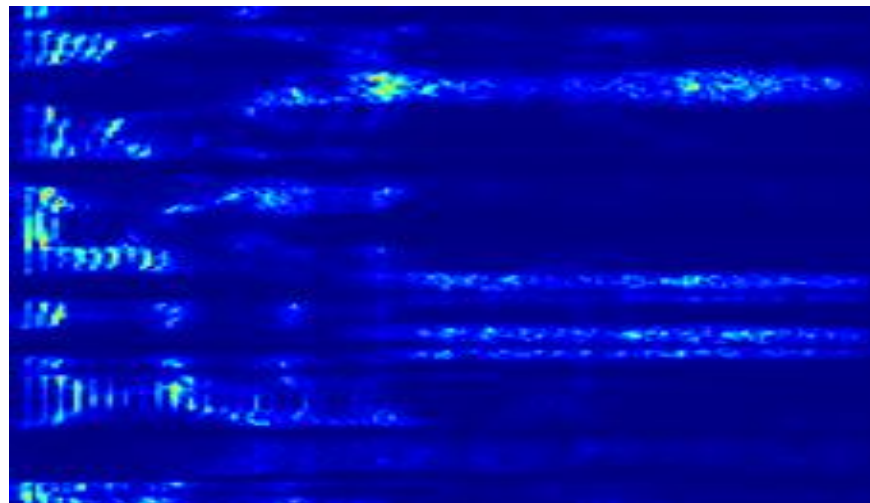# DAY 4: APPLICATIONS OF MACHINE VISION
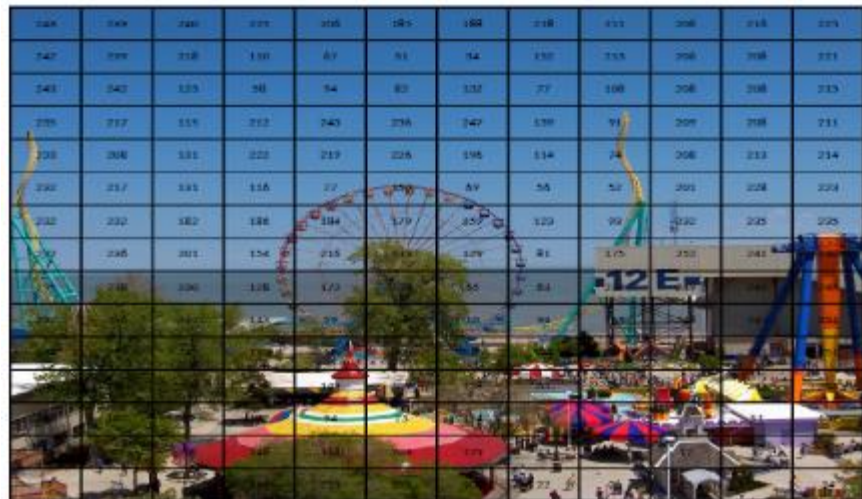
## KE5108: DEVELOPING INTELLIGENT SYSTEMS FOR PERFORMING BUSINESS ANALYTICS

Dr TIAN Jing

tianjing@nus.edu.sg
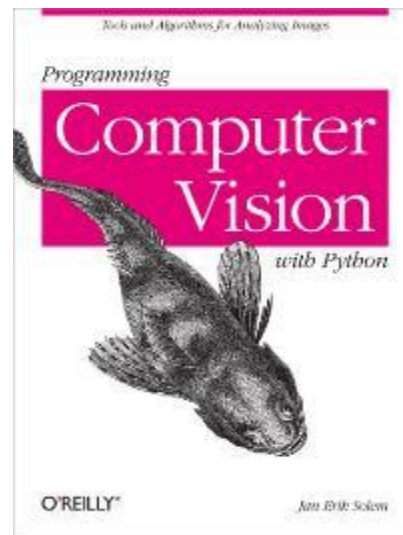
# Reference

- R. C. Gonzalez and R. E. Woods, ***Digital Image Processing***, http://www.imageprocessingplace.com/

- **Computer Vision Crash Course**,   Jia-Bin Huang, https://filebox.ece.vt.edu/~jbhuang/

- **Computer Vision: Algorithms and Applications**, Richard Szeliski, http://szeliski.org/Book/

- Programming Computer Vision with Python, http://programmingcomputervision.com/
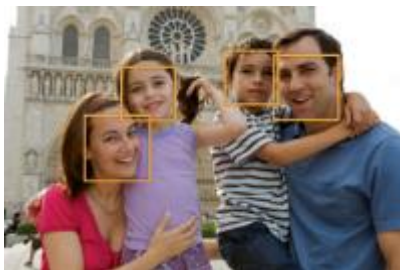
# Outline

- Introduction

- Feature representation: Motion

- Feature representation: Frequency-domain

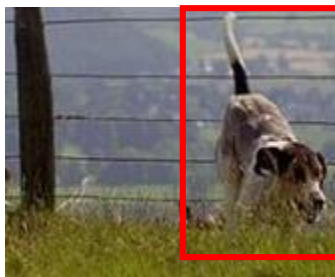- Classification and object detection

# INTRODUCTION

# Computer vision tasks

Face Detection/Recognition

Object Detection

Sports

Object Tracking

Human Pose

Autonomous vehicle

Multi-view Geometry

3D Scene

Vision for Robots

# Outline

| Low Level Task | Mid Level Task | High Level Task |
|---|---|---|
| **Input:** Image<br>**Output:** Image<br><br>**Examples:**<br><br>Noise removal, image sharpening | **Input:** Image<br>**Output:** Attributes<br><br>**Examples:**<br><br>Object recognition, segmentation | **Input:** Attributes/Image<br>**Output:** Understanding<br><br>**Examples:**<br><br>Scene understanding, autonomous navigation |

Image restoration → Image processing and analysis

Image enhancement

Image acquisition

Application domain

Illumination (energy) source

Scene element

Imaging system

(Internal) image plane

Output (digitized) image

Feature representation & description

Object detection and recognition

Scene understanding

# Typical image analytics pipeline

Application domain → Image acquisition → Image enhancement → Image restoration → Image processing and analysis → Feature representation & description → Object detection and recognition → Scene understanding

# Typical image analytics pipeline

Application domain → Image acquisition → Image enhancement → Image restoration → Image processing and analysis → Feature representation & description → Object detection and recognition → Scene understanding

# Typical image analytics pipeline

```
Application domain
  ↓
Image acquisition
  ↓
Image enhancement
  →
Image restoration
  →
Image processing and analysis
  →
Feature representation & description
  ↓
Object detection and recognition
  ↓
Scene understanding
```

$A \ominus B$

$(A \ominus B) \oplus B = A \circ B$

$(A \circ B) \oplus B$ 　　$[(A \circ B) \oplus B] \ominus B = (A \circ B) \cdot B$

# Typical image analytics pipeline



**Image restoration** → **Image processing and analysis**

**Image enhancement**

**Image acquisition**

Application domain

**Feature representation & description**

**Object detection and recognition**

**Scene understanding**

# Typical image analytics pipeline



Application domain → Image acquisition → Image enhancement → Image restoration → Image processing and analysis → Feature representation & description → **Object detection and recognition** → Scene understanding

Labels in figure: bouquet of red flowers, tablet, bottle of water, glass of water with ice and lemon, cup of coffee, dining table with breakfast items, plate of fruit, banana slices, fork, a person sitting at a table

Source: https://adeshpande3.github.io/assets/Caption.png

```
                    ┌──────────────┐      ┌──────────────┐
                    │    Image     │ ───▶ │    Image     │
                    │  restoration │      │  processing  │ ──┐
                    │              │      │ and analysis │   │
                    └──────────────┘      └──────────────┘   │
                           ▲                                  ▼
  ┌──────────────┐         │                          ┌──────────────┐
  │    Image     │ ────────┘                          │   Feature    │
  │ enhancement  │                                    │representation│
  │              │                                    │& description │
  └──────────────┘                                    └──────────────┘
         ▲                                                    │
         │                                                    ▼
  ┌──────────────┐                                    ┌──────────────┐
  │    Image     │                                    │   Object     │
  │ acquisition  │                                    │ detection and│
  │              │                                    │ recognition  │
  └──────────────┘                                    └──────────────┘
         ▲                                                    │
         │                                                    ▼
  Application domain                                  ┌──────────────┐
                                                      │    Scene     │
                                                      │understanding │
                                                      └──────────────┘
```

"A Tabby cat is leaning on a wooden table, with one paw on a laser mouse and the other on a black laptop"

# FEATURE: FREQUENCY-DOMAIN

Weighted summation

$$\left( \begin{array}{ccc} 255 \times 0.0625 & + & 255 \times 0.125 & + & 255 \times 0.0625 \\ + & 209 \times 0.125 & + & 228 \times 0.25 & + & 155 \times 0.125 \\ + & 84 \times 0.0625 & + & 61 \times 0.125 & + & 35 \times 0.0625 \end{array} \right)$$

$$= 181$$

kernel:

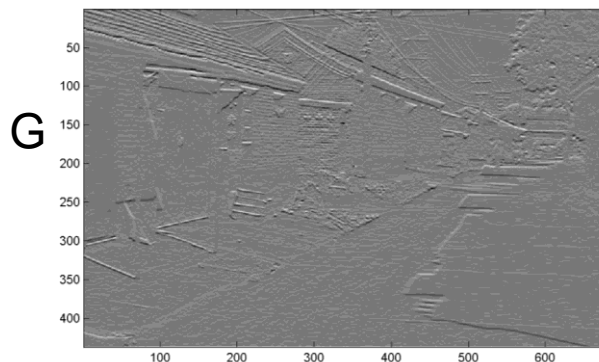blur ▼

input image

output image

Source: http://setosa.io/ev/image-kernels/

* Is the filtering operator

Example question:

Fill in the blanks

```
a)  _  =  D  *  B
b)  F  =  D  *  _
```

A

B

C

D

E

F

G

H

I

# What do you see?

Hybrid Image                Low-passed Image    ➕    High-passed Image

# Inspiration from human perception

- Early processing in humans perception filters for orientations and scales of frequency.



Early visual processing: Multi-scale edge and blob filters

# Convolution theorem

Example: $f$ and $g$ are functions defined in spatial domain, while $F$ and G are their corresponding functions defined in Fourier domain

$$f(x,y) * g(x,y) \Leftrightarrow F(u,v)G(u,v)$$

In words: the Fourier transform of the convolution of two functions is the product of their individual Fourier transforms

Because linear filtering operations can be carried out by simple multiplications in the Fourier domain

# Convolution theorem



$f(x,y)$

Gaussian scale=3 pixels

$g(x,y)$

Fourier transform

Inverse Fourier transform

$|F(u,v)|$

$|G(u,v)|$

Gaussian 

Box filter 



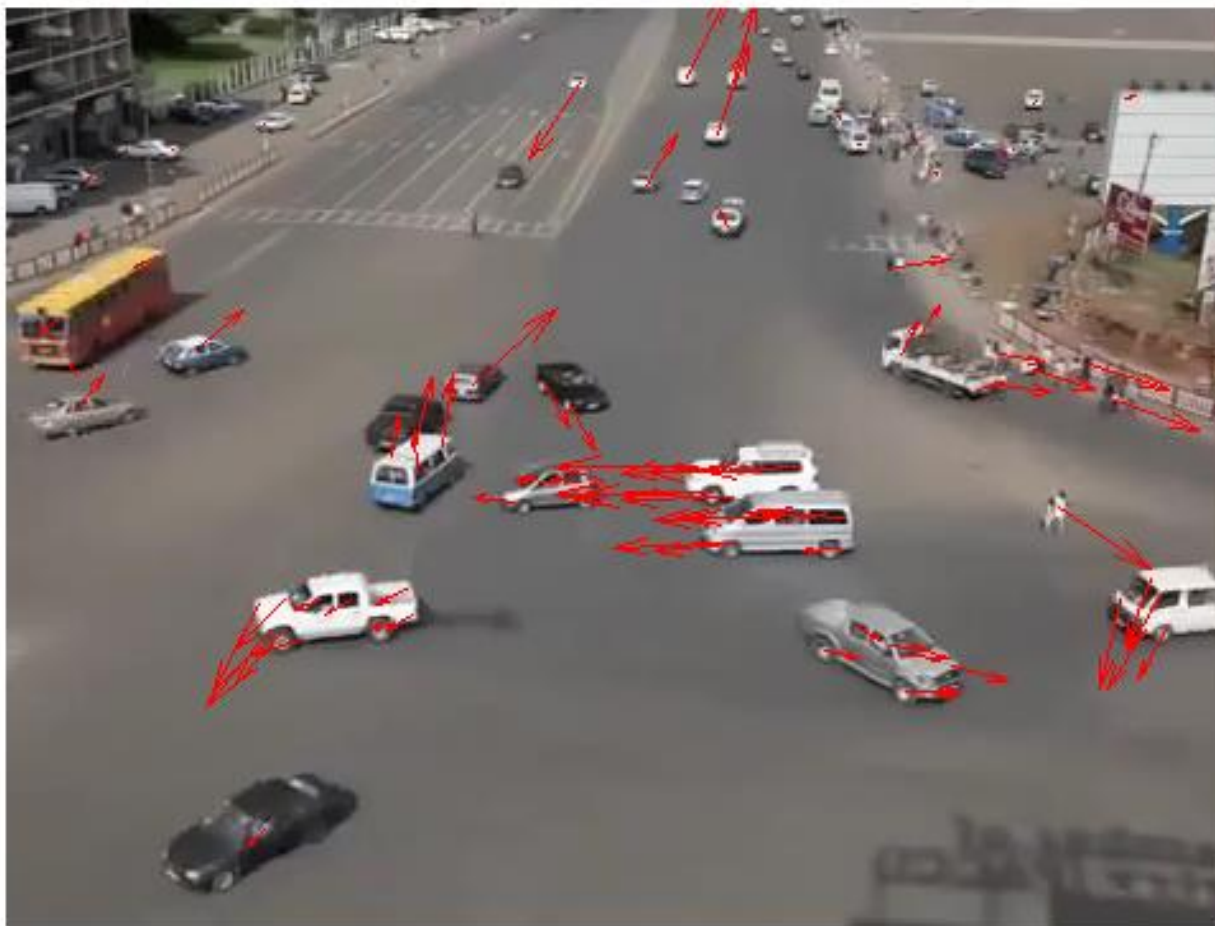| Box filter (spatial) | Frequency domain magnitude | Gaussian filter (spatial) | Frequency domain magnitude |
|---|---|---|---|

# FEATURE: MOTION

- A video is a sequence of frames captured over time

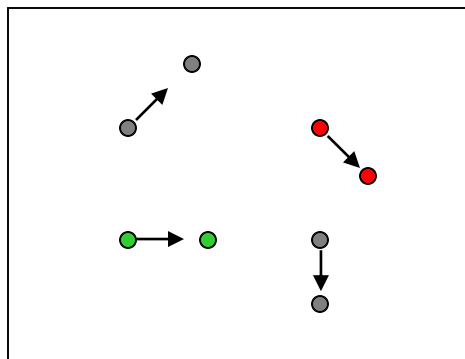- Now our image data is a function of space (x, y) and time (t)
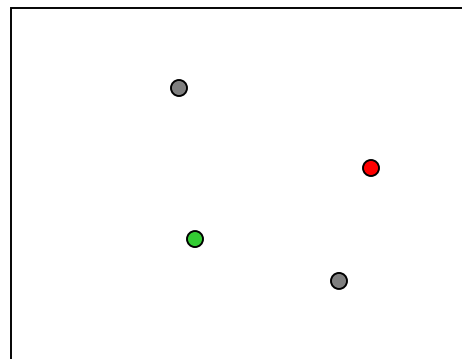


$I(x,y,t)$

$t$

# Optical flow

- Definition: optical flow is the **apparent** motion of brightness patterns in the image

- Note: apparent motion can be caused by lighting changes without any actual motion

**GOAL:** Estimate image motion at each pixel from optical flow.
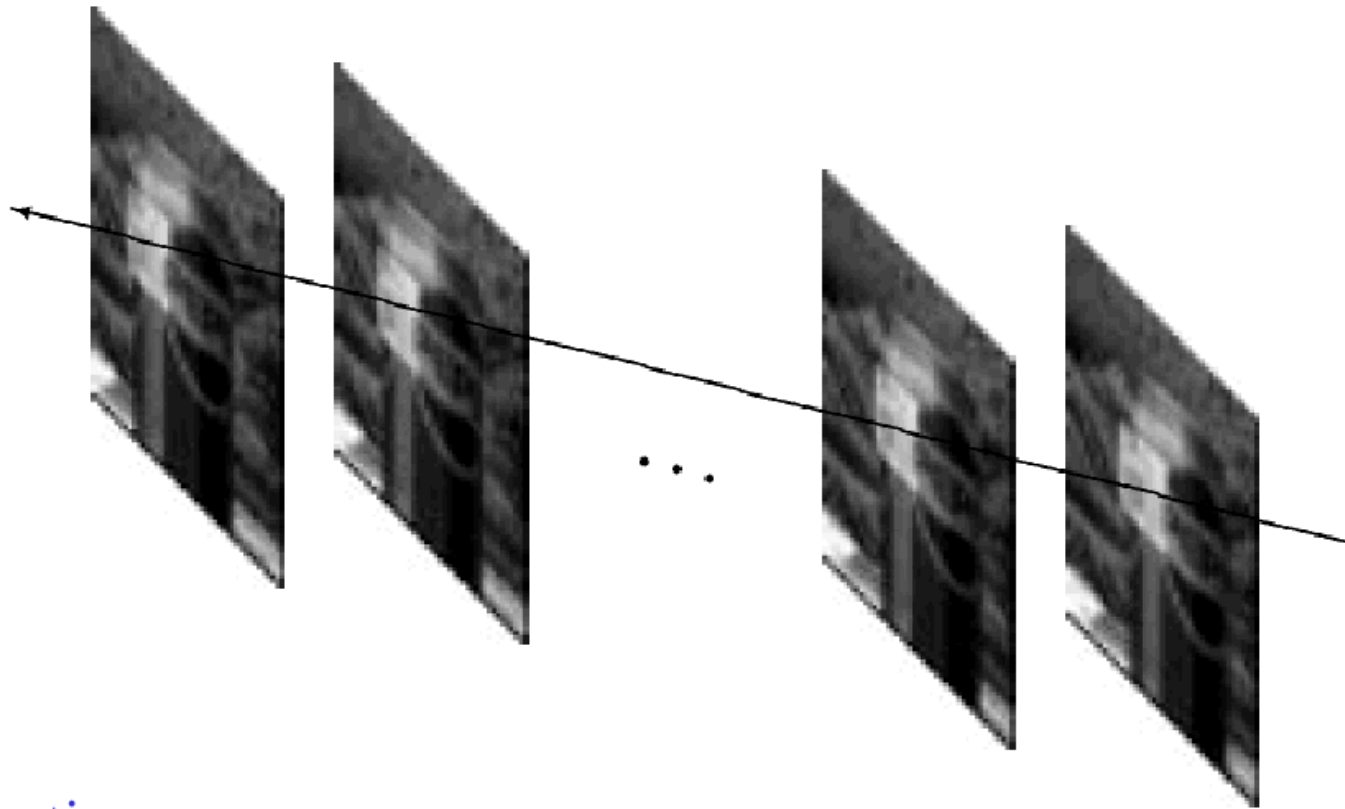
# **Estimating optical flow**



$$I(x,y,t-1) \qquad\qquad I(x,y,t)$$

- Given two subsequent frames, estimate the apparent motion vector field $u(x,y)$, $v(x,y)$ between them

- Key assumptions
  - **Brightness constancy:**  projection of the same point looks the same in every frame
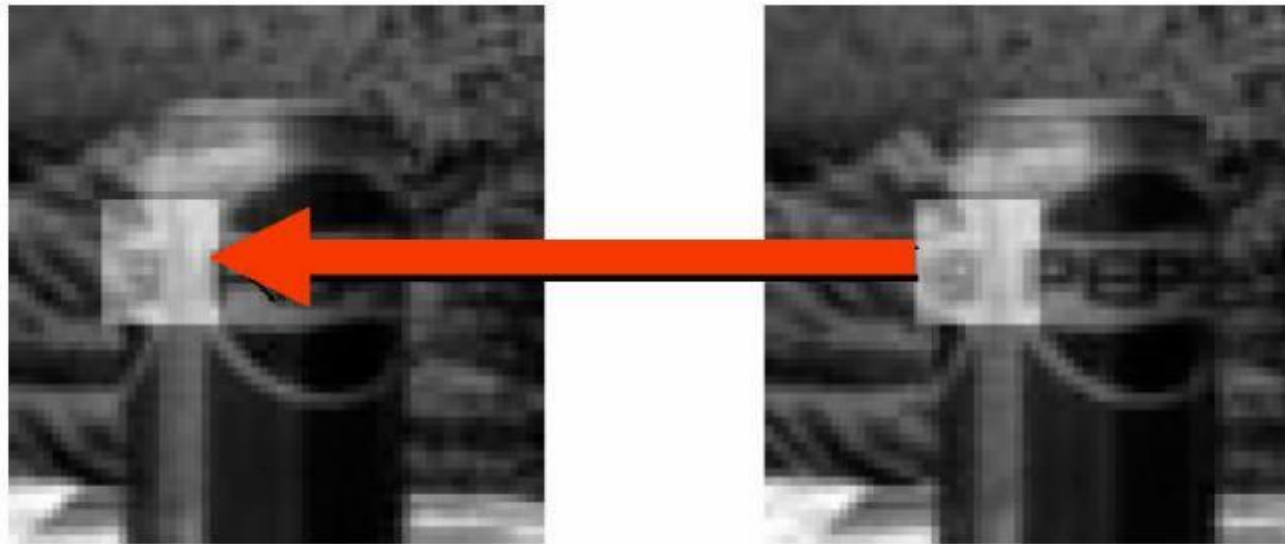  - **Small motion:**  points do not move very far; the length of the vector $u(x,y)$, $v(x,y)$ are small.

# Small motion



Assumption:
 The image motion of a surface patch changes gradually over time.

# Brightness constancy



## Assumption

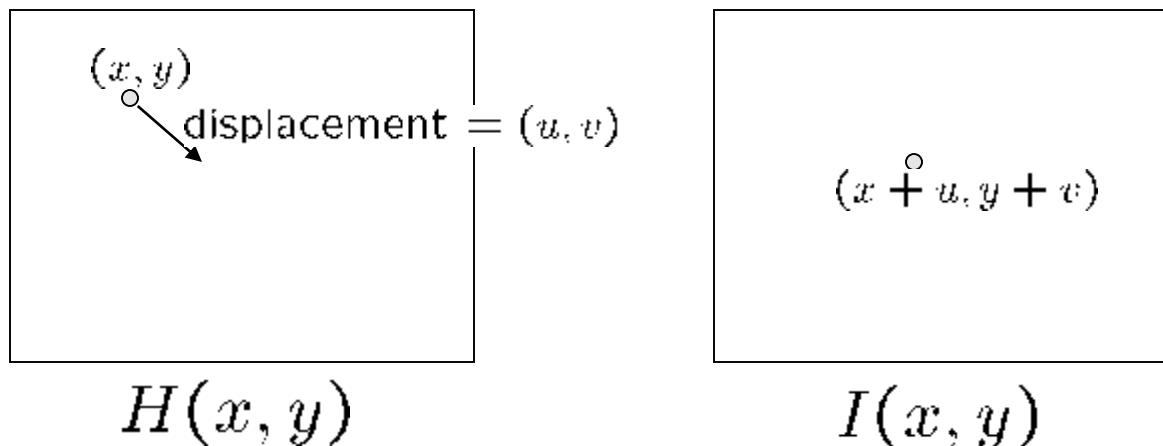Image measurements (e.g. brightness) in a small region remain the same although their location may change.

$$I(x+u, y+v, t+1) = I(x, y, t)$$

(assumption)

$$H(x, y)$$

$$I(x, y)$$

- brightness constancy:   H(x,y)=I(x+u, y+v)

- small motion: suppose we take the Taylor series expansion of I:

$$I(x+u, y+v) = I(x, y) + \frac{\partial I}{\partial x}u + \frac{\partial I}{\partial y}v + \boxed{\text{higher order terms}}$$

$$\approx I(x, y) + \frac{\partial I}{\partial x}u + \frac{\partial I}{\partial y}v$$

Combining these two equations

shorthand: $I_x = \frac{\partial I}{\partial x}$

$$0 = I(x + u, y + v) - H(x, y)$$

$$\approx I(x, y) + I_x u + I_y v - H(x, y)$$

$$\approx (I(x, y) - H(x, y)) + I_x u + I_y v$$

$$\approx I_t + I_x u + I_y v$$

$$\approx I_t + \nabla I \cdot [u \ v]$$

In the limit as u and v go to zero, this becomes exact

$$0 = I_t + \nabla I \cdot [\tfrac{\partial x}{\partial t} \ \tfrac{\partial y}{\partial t}]$$

- How to get more equations for a pixel?

- **Spatial coherence constraint:**

- Assume the pixel's neighbors have the same (u,v)

  - If we use a 5x5 window, that gives us 25 equations per pixel

$$0 = I_t(\mathbf{p_i}) + \nabla I(\mathbf{p_i}) \cdot [u \ v]$$

$$\begin{bmatrix} I_x(\mathbf{p_1}) & I_y(\mathbf{p_1}) \\ I_x(\mathbf{p_2}) & I_y(\mathbf{p_2}) \\ \vdots & \vdots \\ I_x(\mathbf{p_{25}}) & I_y(\mathbf{p_{25}}) \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = - \begin{bmatrix} I_t(\mathbf{p_1}) \\ I_t(\mathbf{p_2}) \\ \vdots \\ I_t(\mathbf{p_{25}}) \end{bmatrix}$$

B. Lucas and T. Kanade. An iterative image registration technique with an application to stereo vision.
In *Proceedings of the International Joint Conference on Artificial Intelligence*, pp. 674–679, 1981.

- Overconstrained linear system

$$\begin{bmatrix} I_x(\mathbf{p_1}) & I_y(\mathbf{p_1}) \\ I_x(\mathbf{p_2}) & I_y(\mathbf{p_2}) \\ \vdots & \vdots \\ I_x(\mathbf{p_{25}}) & I_y(\mathbf{p_{25}}) \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = - \begin{bmatrix} I_t(\mathbf{p_1}) \\ I_t(\mathbf{p_2}) \\ \vdots \\ I_t(\mathbf{p_{25}}) \end{bmatrix}$$

$$A \quad d = b$$
$$\text{25x2} \quad \text{2x1} \quad \text{25x1}$$

Least squares solution for *d* given by $(A^T A)\, d = A^T b$

$$\underbrace{\begin{bmatrix} \sum I_x I_x & \sum I_x I_y \\ \sum I_x I_y & \sum I_y I_y \end{bmatrix}}_{A^T A} \begin{bmatrix} u \\ v \end{bmatrix} = - \underbrace{\begin{bmatrix} \sum I_x I_t \\ \sum I_y I_t \end{bmatrix}}_{A^T b}$$

# Iterative refinement at same scale

- Iterative Lukas-Kanade Algorithm
    1. Estimate velocity at each pixel by solving Lucas-Kanade equations
    2. Warp I(t-1) towards I(t) using the estimated flow field
        - *use image warping techniques*
    3. Repeat until convergence

run iterative
Lukas-Kanade
algorithm

warp & upsample

run iterative
Lukas-Kanade
algorithm

image A

image B

Gaussian pyramid of image A at (t-1)    Gaussian pyramid of image B at (t)

- In other words, in what situations does the displacement of pixel patches not represent physical movement of points in space?

- A uniform rotating
  - nothing seems to move, yet it is rotating

- Changing directions or intensities of lighting can make things seem to move
  - for example, if the specular highlight on a rotating sphere moves.

# MACHINE LEARNING

# Machine learning

**Objective: Looking for a function!**

- Speech Recognition

$$f\left( \begin{array}{c} \text{[waveform]} \end{array} \right) = \text{"How are you"}$$

- Image Recognition

$$f\left( \begin{array}{c} \text{[cat image]} \end{array} \right) = \text{"Cat"}$$

- Playing Go

$$f\left( \begin{array}{c} \text{[Go board]} \end{array} \right) = \text{"5-5"} \quad \text{(next move)}$$

- Chatbot

$$f\left( \quad \text{"Hi"} \quad \right) = \quad \text{"Hello"}$$

(what the user said)   (system response)

- Object **Classification**
  what object ?

- Object **Detection**
  object or no-object ?

- Instance **Recognition** ?
  who (or what) is it ?

- **Sub-category** analysis
  which object type ?

- Sequence { **Recognition | Classification** } ?
  what is happening / occurring ?

http://pascallin.ecs.soton.ac.uk/challenges/VOC/

{people | vehicle | … intruder ….}

{face | vehicle plate| gait …. → biometrics}

GN05 X·JA

{gender | type | species | age ……}

# Image classification



Training

Training Images → Image Features → Classifier Training → Trained Classifier

Training Labels → Classifier Training

Testing

Test Image → Image Features → Trained Classifier → Prediction result Outdoor

1.  Select / develop features:  SURF, HoG, SIFT, …
2.  Add on top of this Machine Learning for multi-class recognition and train classifier



| | Feature Extraction: SIFT, HoG... | → | Detection, Classification Recognition | → | |

Classical computer vision feature definition is domain-specific and time-consuming

# What are the right features?



Image gradients → Keypoint descriptor

Point sampling strategy — Color descriptor computation — Bag-of-words model

Image — Harris-Laplace detector — Histograms, ColorSIFT, ... — Vector quantisation

# What are the right features?

- Object: shape
  - Local shape info, shading, shadows, texture
- Scene: geometric layout
  - linear perspective, gradients, line segments
- Material properties: Color, texture
- Action: motion
  - Optical flow, tracked points

# Histogram

- Histogram of an image provides the frequency of the brightness (intensity) value in the image.

```
def histogram(im):
    h = np.zeros(255)
    for row in im.shape[0]:
        for col in im.shape[1]:
            val = im[row, col]
            h[val] += 1
```



| 0 | 256 |
| --- | --- |

Count: 10192  Min: 9
Mean: 133.711  Max: 255
StdDev: 55.391  Mode: 178 (180)

Count: 10192  Min: 11
Mean: 104.637  Max: 254
StdDev: 89.862  Mode: 23 (440)

- Color



L*a*b* color space

HSV color space

- Texture (filter banks)

# Histogram

## Gradients

## "Bag of visual words"



Image gradients

Keypoint descriptor
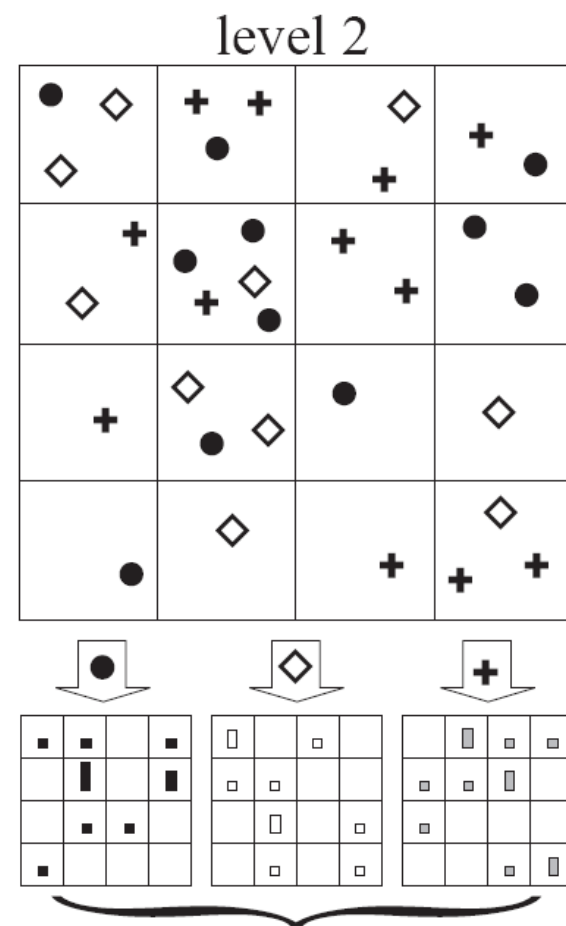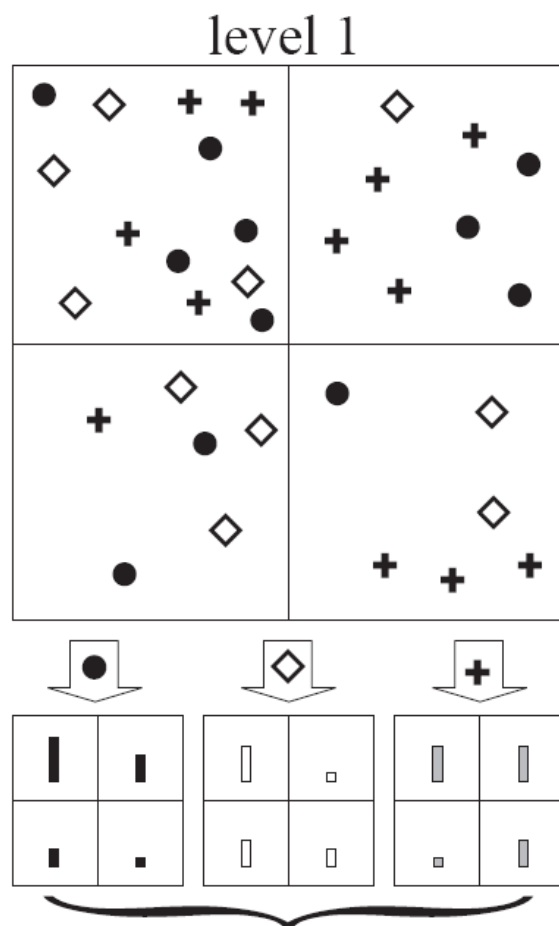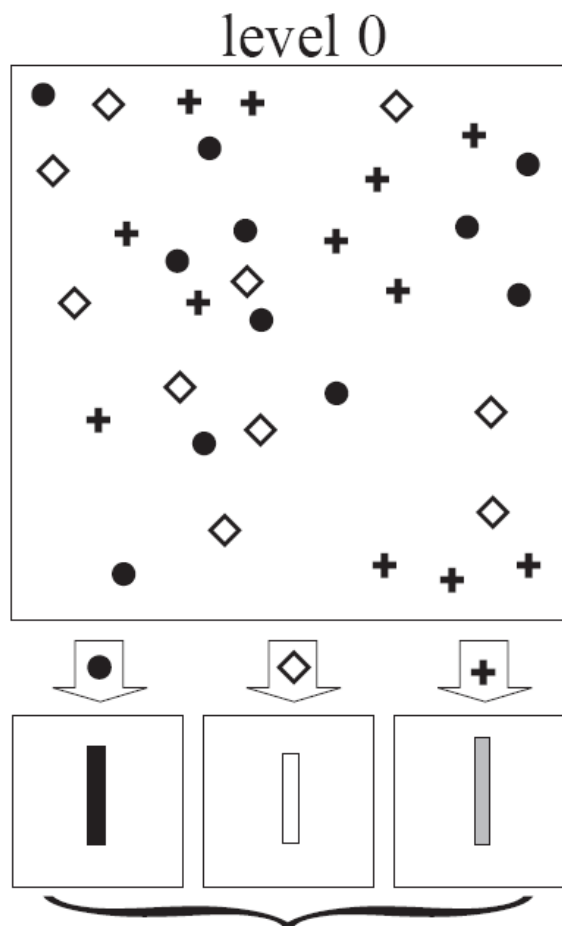
All of these images have the same color histogram

Compute histogram in each spatial bin

# Spatial pyramid

# OBJECT DETECTION
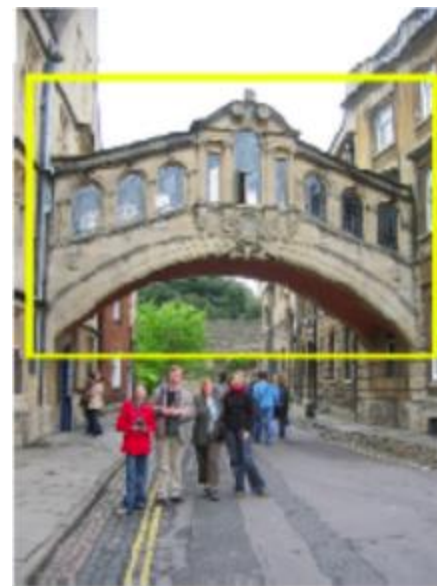
# Category vs. instance recognition

Category:
- Find all the people
- Find all the buildings
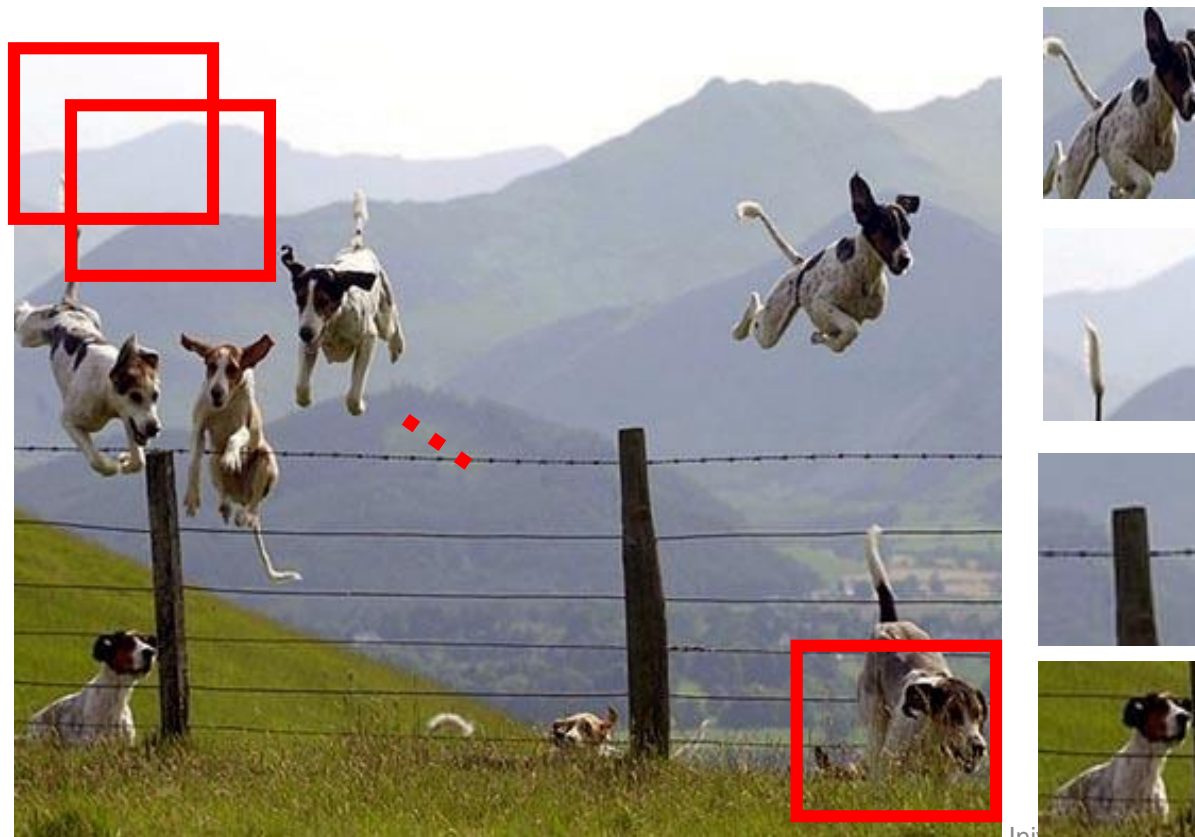- Often within a single image
- Often 'sliding window'



Instance:
- Is this face James?
- Find this specific famous building
- Often within a database of images

# Object category detection

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



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

# Challenges in object detection
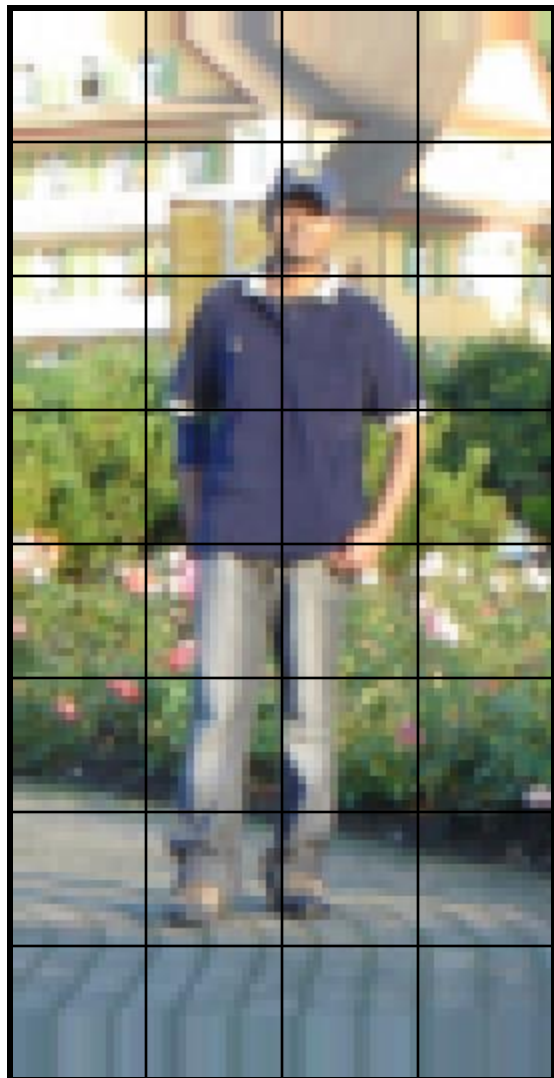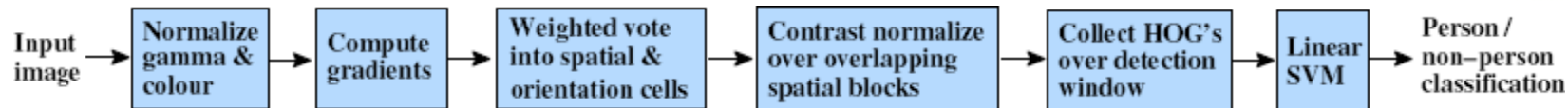
Illumination

Object pose

'Clutter'

Occlusions

Intra-class appearance

Viewpoint

Input image → Normalize gamma & colour → Compute gradients → Weighted vote into spatial & orientation cells → Contrast normalize over overlapping spatial blocks → Collect HOG's over detection window → Linear SVM → Person / non–person classification

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

Input image → Normalize gamma & colour → Compute gradients → Weighted vote into spatial & orientation cells → Contrast normalize over overlapping spatial blocks → Collect HOG's over detection window → Linear SVM → Person / non–person classification
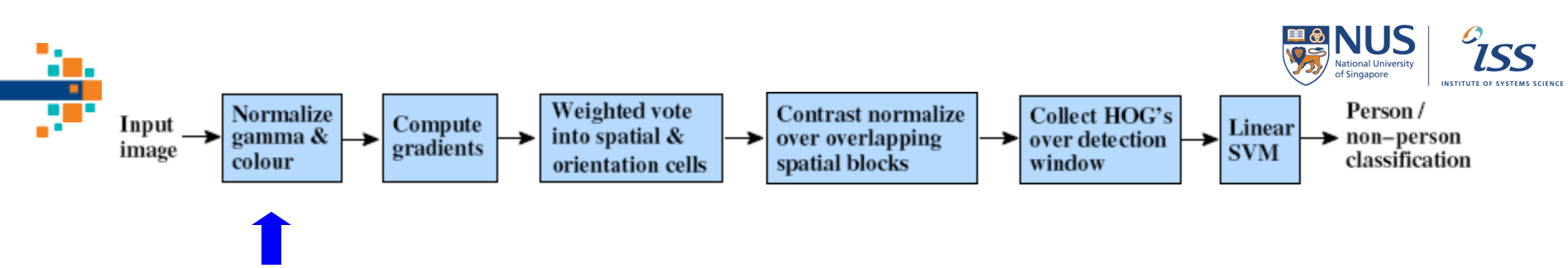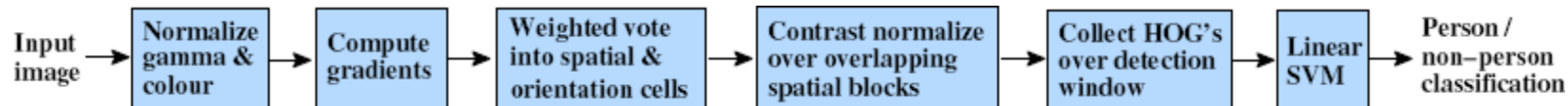
- Tested with
  - RGB
  - LAB
  - Grayscale

  Slightly better performance vs. grayscale

- Gamma Normalization and Compression
  - Square root
  - Log

  Very slightly better performance vs. no adjustment

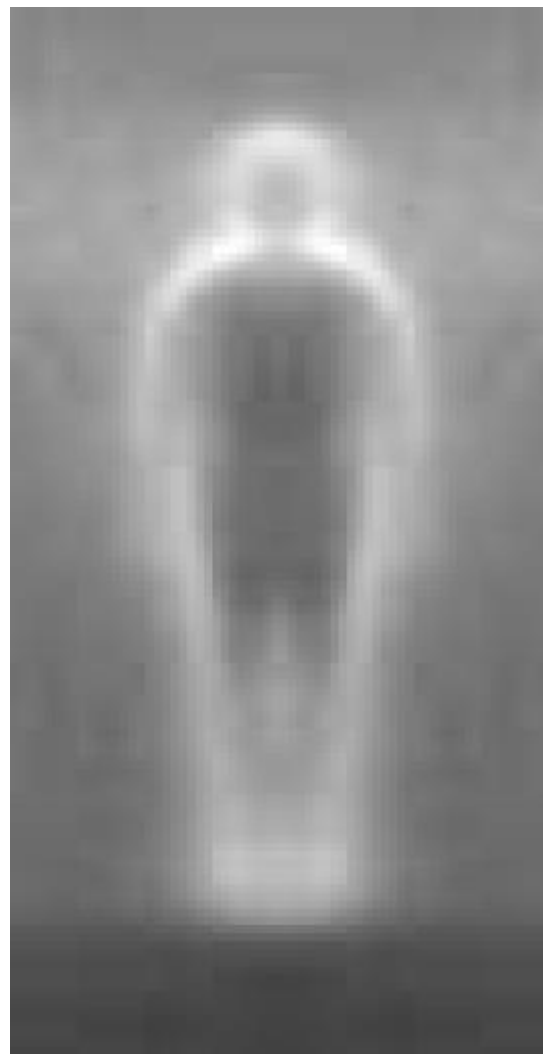Input image → Normalize gamma & colour → Compute gradients → Weighted vote into spatial & orientation cells → Contrast normalize over overlapping spatial blocks → Collect HOG's over detection window → Linear SVM → Person / non–person classification

| -1 | 0 | 1 |

centered

| -1 | 1 |

uncentered

| 1 | -8 | 0 | 8 | -1 |

cubic-corrected

| 0 | 1 |
| -1 | 0 |

diagonal

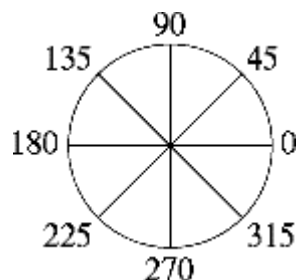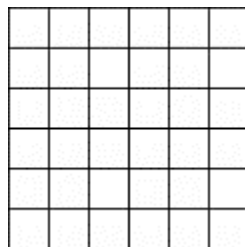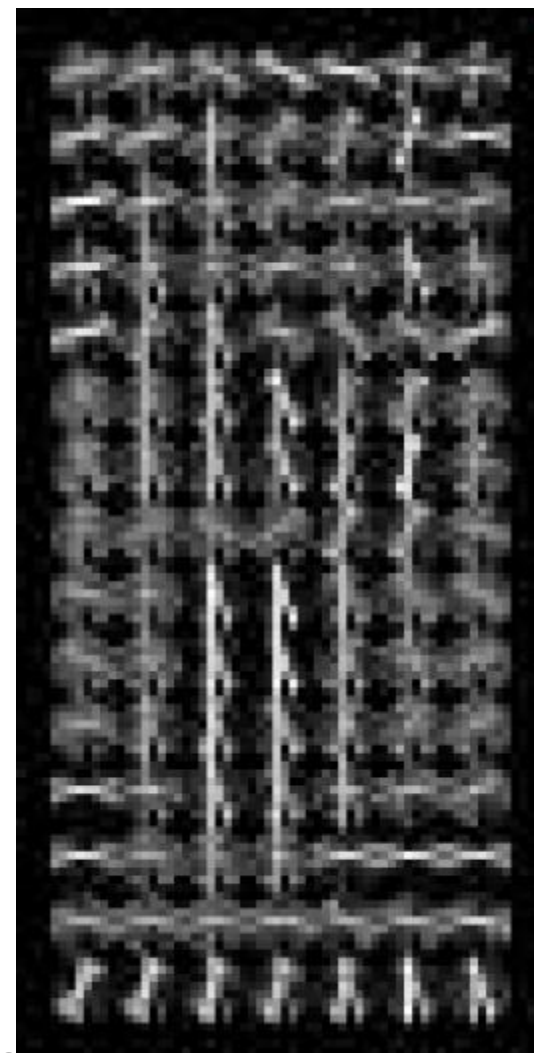| -1 | 0 | 1 |
| -2 | 0 | 2 |
| -1 | 0 | 1 |

Sobel

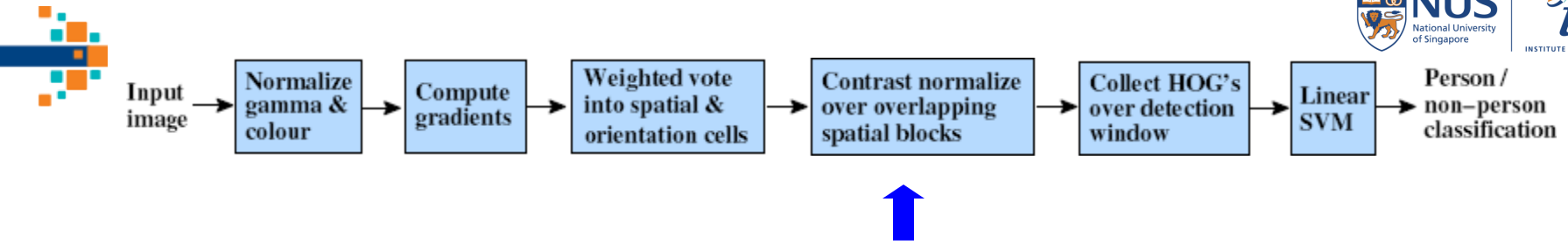- # Histogram of gradient orientations

Orientation: 9 bins
(for unsigned angles)

Histograms in
8x8 pixel cells

- # Votes weighted by magnitude

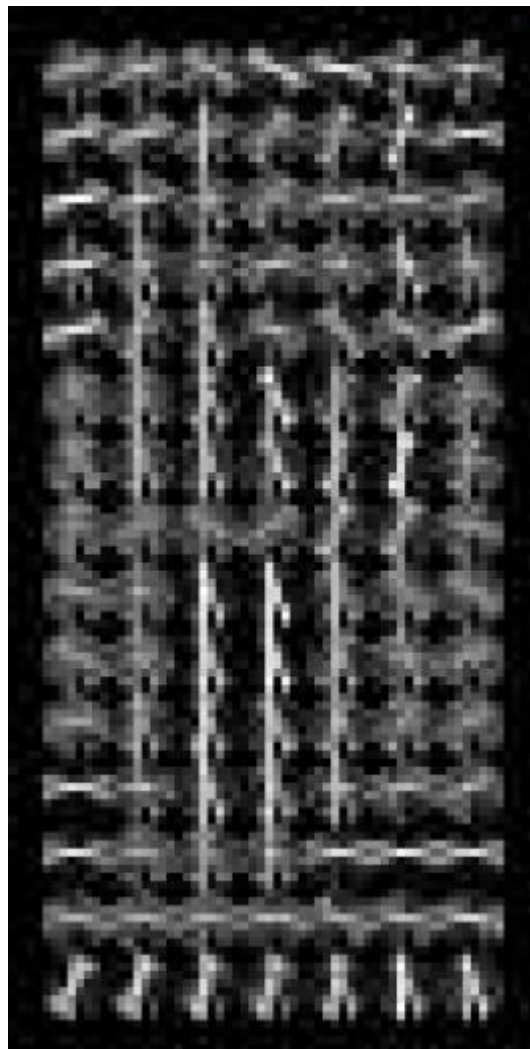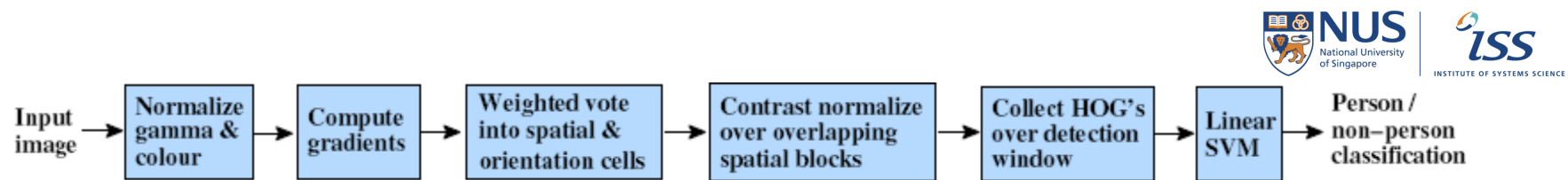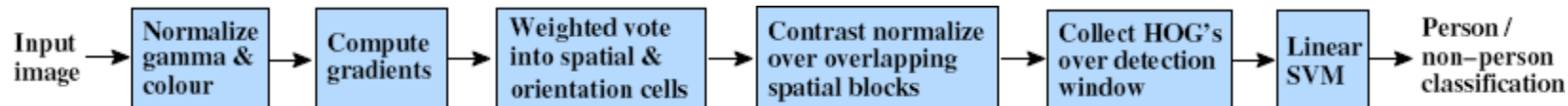| Input image | → | Normalize gamma & colour | → | Compute gradients | → | Weighted vote into spatial & orientation cells | → | Contrast normalize over overlapping spatial blocks | → | Collect HOG's over detection window | → | Linear SVM | → | Person / non–person classification |

## R-HOG

Cell

Block

Normalize with respect to surrounding cells

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

Input image → Normalize gamma & colour → Compute gradients → Weighted vote into spatial & orientation cells → Contrast normalize over overlapping spatial blocks → Collect HOG's over detection window → Linear SVM → Person / non–person classification

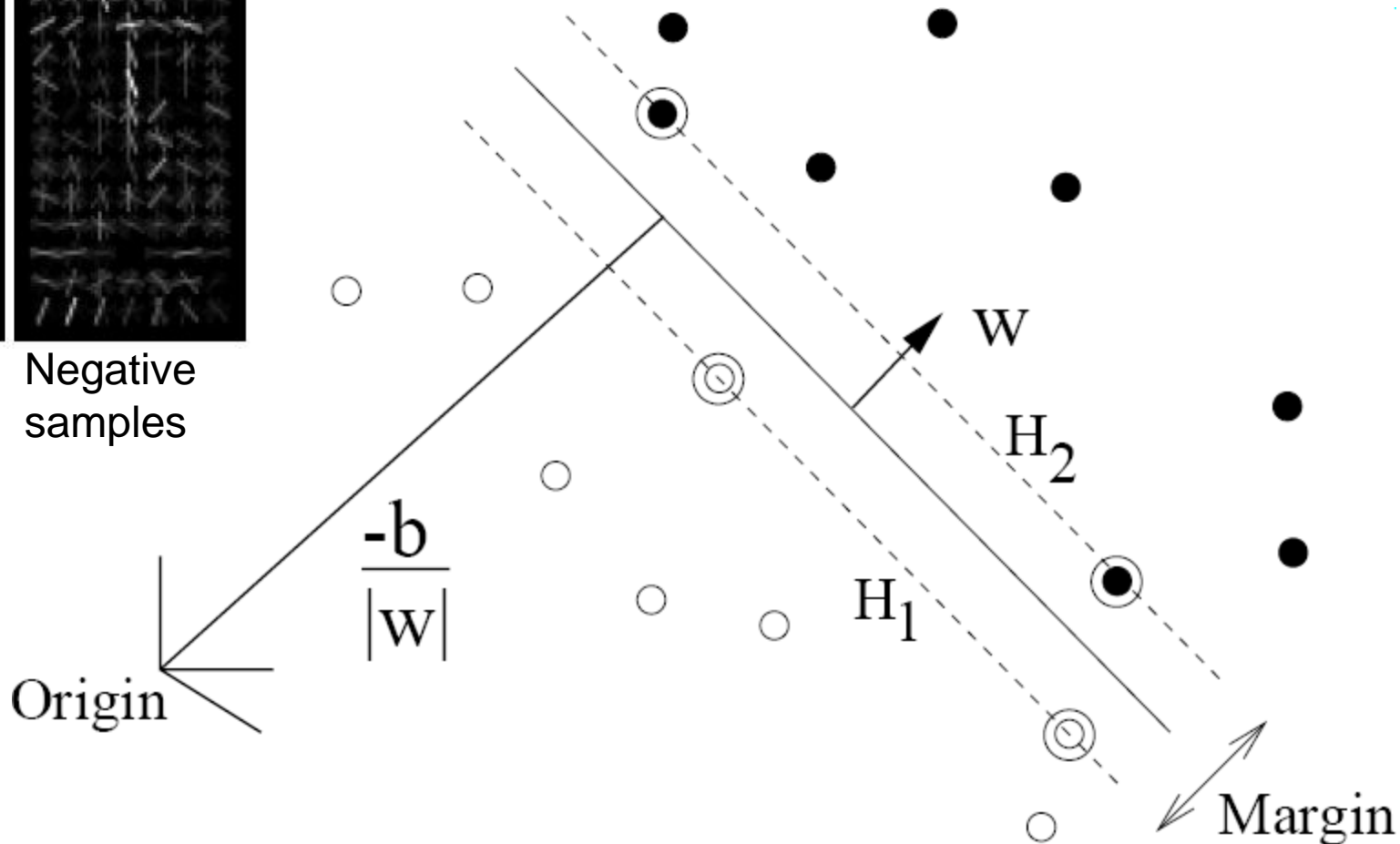Positive samples    Negative samples

$$\frac{-b}{|w|}$$

Origin

W

$H_2$

$H_1$

Margin
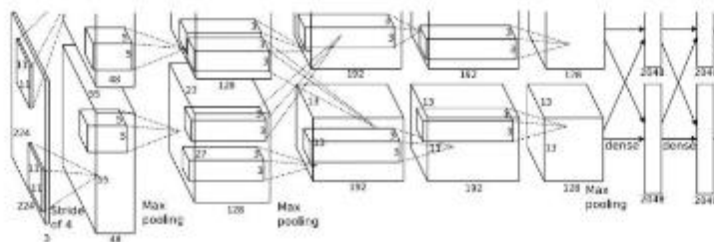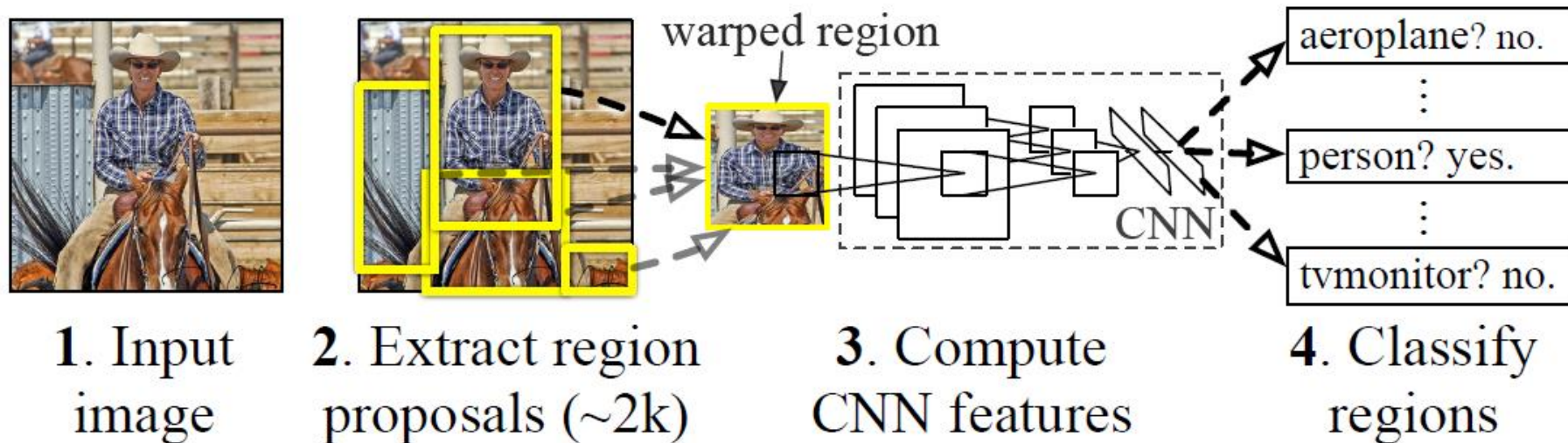
# CNN as feature extractor



Apply a CNN to many different crops of the image, CNN classifies each crop as object or background

Dog? NO
Cat? NO
Background? YES

# CNN as feature extractor



1. Input image
2. Extract region proposals (~2k)
3. Compute CNN features
4. Classify regions
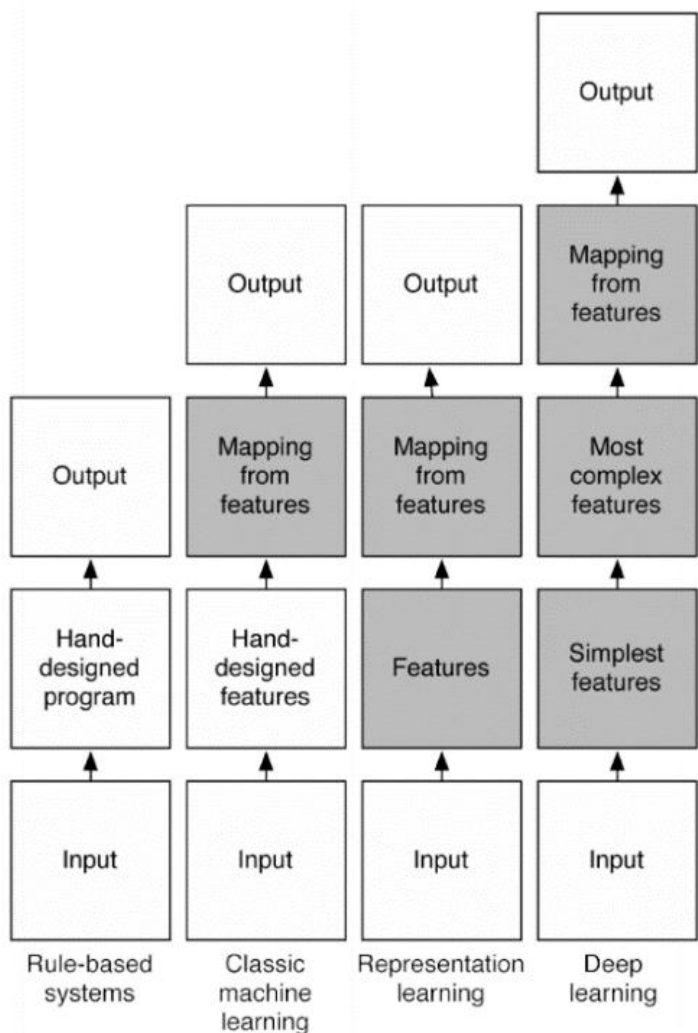
- Replace sliding windows with "selective search" region proposals
- Extract rectangles around regions
- Extract features with fine-tuned CNN (that was initialized with network trained on ImageNet before training)
- Classify last layer of network features with SVM, refine bounding box localization (bbox regression) simultaneously

http://arxiv.org/pdf/1311.2524.pdf

"Deep learning allows computational models that are composed of **multiple processing layers to learn representations of data with multiple levels of abstraction**."
– Yann LeCun, Yoshua Bengio and Geoff Hinton

Y. LeCun, Y. Bengio, G. Hinton, **"Deep Learning"**, Nature, Vol. 521, 28 May 2015

"having had countless ConvNet papers rejected, published and ignored, and occasionally paid attention to, for over 15 years"
-- Yann Lecun

# Summary

- Introduction

- Feature representation: Motion

- Feature representation: Frequency-domain

- Classification and object detection

# Thank You!

Dr TIAN Jing
Email: tianjing@nus.edu.sg