



# La visione

Prof. Alberto Borghese  
Dipartimento di Scienze dell'Informazione  
[borgnese@dsi.unimi.it](mailto:borgnese@dsi.unimi.it)

Università degli Studi di Milano

Slide in parte tratte da: <http://www.andrew.cmu.edu/course/15-491>



# Sommario

- **La visione**
- Le immagini digitali
- Il modello geometrico di una camera
- Segmentazione real-time.



## Le funzioni della visione



A.A. 2008-2009

3/80

<http://homes.dsi.unimi.it/~borghese>



## Computer Vision



**Obiettivo:** determinazione delle proprietà **geometriche**, **fisiche** e **dinamiche** del mondo che ci circonda mediante elaborazione di immagini o sequenze di immagini.

- *Low level vision (o early vision):* estrazione dalle immagini o sequenze di immagini delle informazioni necessarie al livello superiore (features = caratteristiche locali) utili alle elaborazioni successive.
- *High level vision:* riconoscimento, associazione di un significato semantico all'atto del vedere, ricostruzione del movimento degli oggetti.

Nel nostro caso, comportamento reattivo, principalmente low-level vision.

A.A. 2008-2009

4/80

<http://homes.dsi.unimi.it/~borghese>



# Processing visivo



- “Low-level vision (early vision)” – Pre-elaborazione delle immagini (estrazione di feature).
- Calcolo del Movimento degli oggetti sull’immagine (optical flow).
- Estrazione del colore.
- Estrazione della profondita’.
- Riconoscimento di tessiture.
- Contorni (edge)
- “Low-level vision (intermediate representations)”.
- Calcolo delle sorgenti di illuminazione e stima dell’albedo e del colore.
- Forme dai contorni (shape from edges).
- Forme da tessitura (shape from texture).
- Forme da ombreggiatura (shape from shading).
- Stereo-matching.
- Determinazione della struttura 3D e del movimento 3D di oggetti da sequenze di immagini monoculari e da sistemi di specchi (Structure from Motion).
- Ricostruzione 3D da stereo di oggetti della scena.
- Ricostruzione di superfici.
- Parametri geometrici del sistema di visione (calibrazione).
- “High-level vision”.
- Interpretazione e movimento (la visione artificiale deriva storicamente dall’Intelligenza Artificiale).



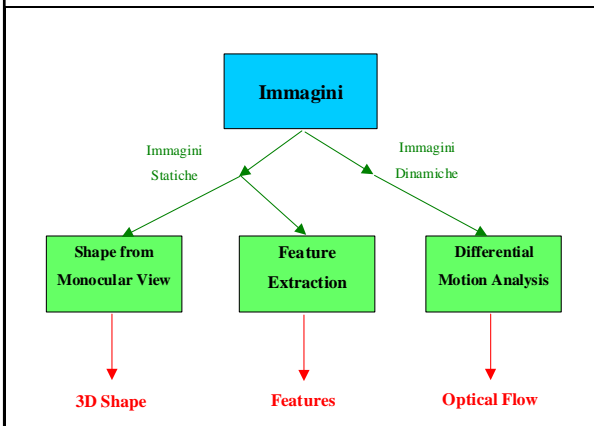
# Visione 3D, Elaborazione di immagini e grafica



**Visione 3D:** Immagine/i => Ricostruzione 3D della scena statica o dinamica ed interpretazione.

• **Grafica 3D:** Modello 3D della scena, statico o dinamico => Visualizzazione.

*Si incontrano sul terreno della visualizzazione 3D.*



L’elaborazione delle immagini costituisce il primo livello di un sistema di visione. Fornisce le features di base.



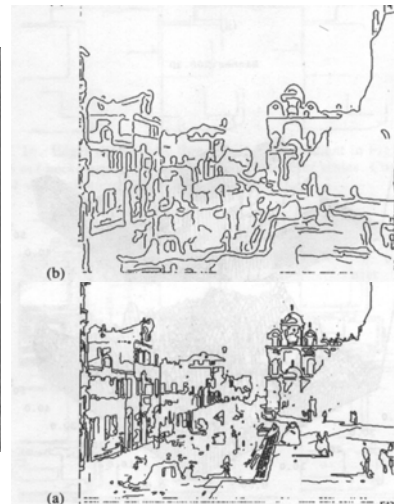
## Cosa sono le features?



- 1) *Località.*
- 2) *Significatività.*
- 3) *Riconoscibilità.*



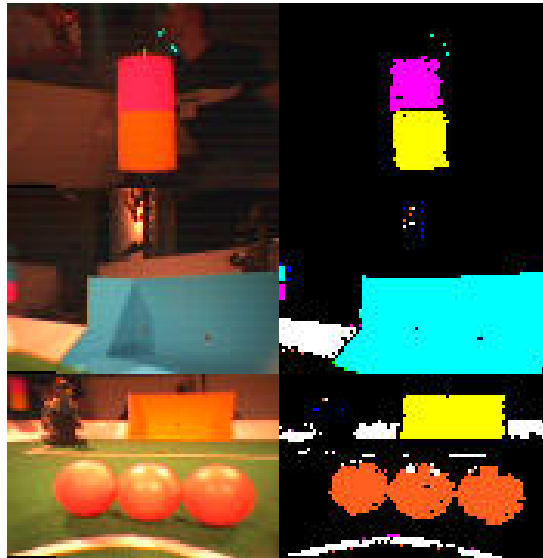
## Riconoscimento dei bordi (edge)



Ricostruzione scena (cf. 3D scanning)



## Estrazione di regioni



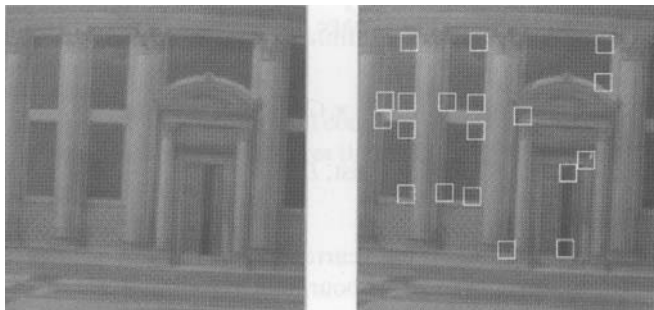
A.A. 2008-2009

9/80

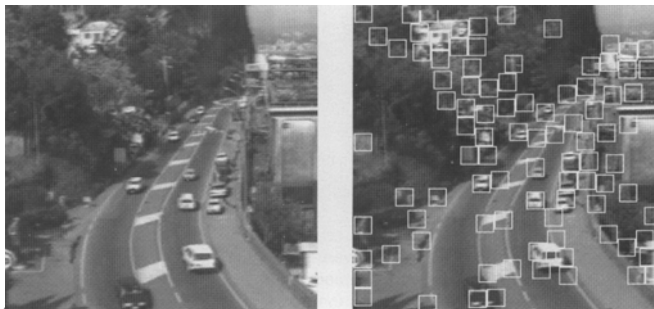
<http://homes.dsi.unimi.it/~borghese>



## Riconoscimento di spigoli



Navigation



A.A. 2008-2009

10/80

<http://homes.dsi.unimi.it/~borghese>



## Il flusso ottico



Navigation

AIBO nella pista

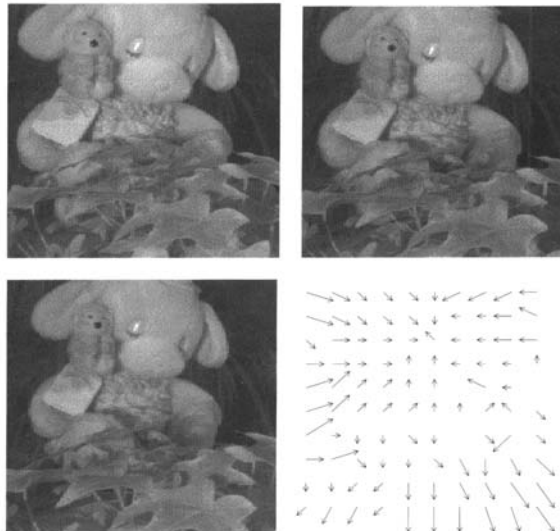


Figure 8.3 Three frames from a long image sequence (left to right and top to bottom) and the optical flow computed from the sequence, showing that the plant in the foreground is moving towards the camera, and the soft toys away from it.

A.A. 2008-2009

12/80

<http://homes.dsi.unimi.it/~borghese>



## I problemi di visione sono mal posti



- Perché non è facile costruire un sistema di visione?

### Difficoltà ad identificare esattamente le feature

- Risoluzione spaziale limitata.
- Gli oggetti reali non sono mai uniformemente illuminati.
- I contorni non sono netti.
- Le superfici non hanno albedo costante.
- L'illuminazione genera campi di irradianza "dificili".

### Difficoltà ad assemblare le feature

### Difficoltà ad interpretare le primitive visive

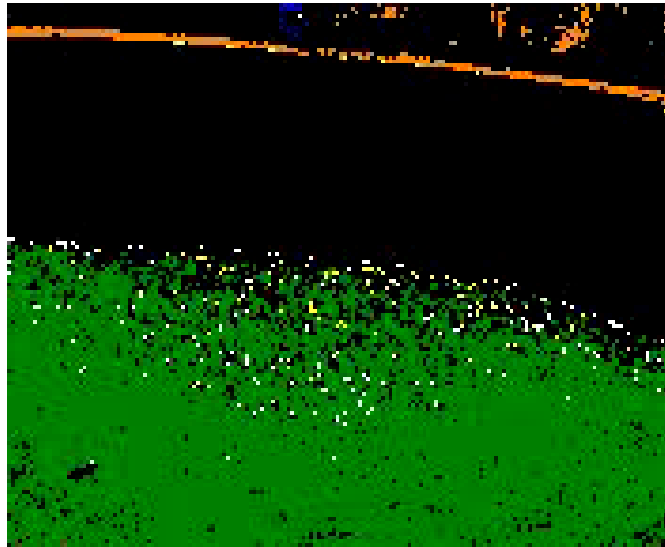
A.A. 2008-2009

13/80

<http://homes.dsi.unimi.it/~borghese>



## La vera Visione



Video



## Sommario



- La visione
- **Le immagini digitali**
- Il modello geometrico di una camera
- Segmentazione real-time.

**IMMAGINI DIGITALI**

Una matematica del pixel non è ancora disponibile.

ny

pixel

nx

IMG = Matrice nx,ny =  $\begin{pmatrix} 142 & 174 & 164 & 144 & \dots & \dots \\ 107 & \dots & & & & \\ \dots & & & & & \\ \dots & & & & & \end{pmatrix}$  ny

A.A. 2008-2009

16/80

**IMMAGINI DIGITALI**

**Quantizzazione:** n bit  $\Rightarrow 2^n$  colori

ES.: 3 bit  $\Rightarrow 2^3 = 8$  'gradini'

2 bit  $\Rightarrow 2^2 = 4$  'gradini'

Colormap

8 bit  $\Rightarrow 256$  colori

Colormap

3 bit  $\Rightarrow 8$  colori

A.A.

17/80

s.dsi.unimi.it/~borghese





## Colore



Colour is the colour which is perceived, seen, that is the colour which is reflected by the objects surface. Color is coded with three parameters (three channels).

Colour images are represented as additive mixture of Red Green Blue (additive mix).

*Another important coding is:*

**Y. Brightness.** Intensità del colore. It can be viewed as the colour of the image in B/W. It is due to the illumination intensity.

**Cb, Cr.** Quantità di blu e di rosso all'interno di un'immagine.

Y – Brightness.

Y

U, V – Color.

Cb, Cr



## La trasformazione tra spazio YCbCr ed RGB



$$Y' = K_r * R' + (1 - K_r - K_b) * G' + K_b * B'$$

$$C_b = 0.5 * (B' - Y') / (1 - K_b)$$

$$C_r = 0.5 * (R' - Y') / (1 - K_r)$$

$$0 \leq Y', R, G, B \leq +1$$

$$-0.5 \leq C_b, C_r \leq +0.5$$

In funzione di  $K_r$  e  $K_b$   
abbiamo diversi  
standard



## Standard Video (ITU-R BT.601) trasformazione diretta



$$K_b = 0.114$$

$$K_r = 0.299$$

$$Y = + 0.299 * R + 0.587 * G + 0.114 * B$$

$$C_b = - 0.168736 * R - 0.331264 * G + 0.5 * B$$

$$C_r = + 0.5 * R - 0.418688 * G - 0.081312 * B$$

$$0 \leq Y', R, G, B \leq +1$$

$$-0.5 \leq C_b, C_r \leq +0.5$$



## Standard Video (ITU-R BT.601) trasformazione inversa



$$R = Y + 1.402 * (C_r - 128)$$

$$G = Y - 0.34414 * (C_b - 128) - 0.71414 * (C_r - 128)$$

$$B = Y + 1.772 * (C_b - 128)$$

$$0 \leq Y, C_b, C_r, R, G, B \leq +255$$

Diversi coefficienti per ITU-R NT.709 - HDTV:

$$K_b = 0.0722$$

$$K_r = 0.2126$$



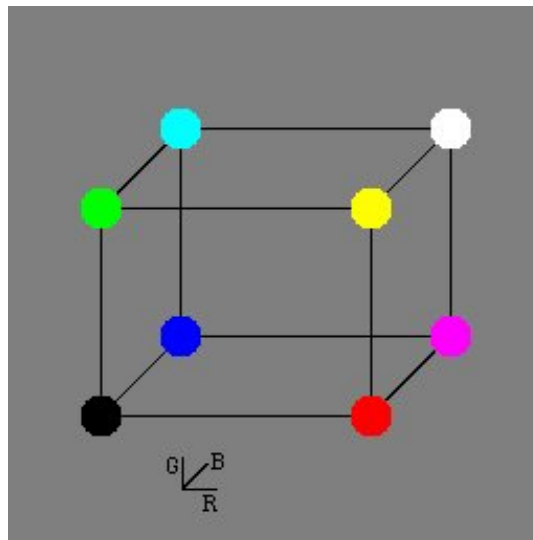
### Colors (examples in RGB)

	White	(R=255, G=255, B=255)
	Light Grey	(R=100, G=100, B=100)
	Dark grey	(R=200, G=200, B=200)
	Black	(R=0, G=0, B=0)
	Red	(R=255, G=0, B=0)
	Yellow	(R=255, G=255, B=0)
	Pale blue	(R=0, G=255, B=255)
	Green	(R=0, G=200, B=0)

A.A. 2008-2009



## Color Spaces - RGB



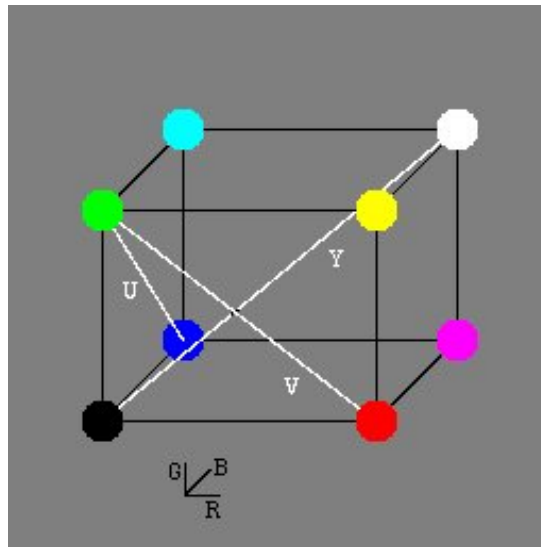
A.A. 2008-2009

24/80

<http://homes.dsi.unimi.it/~borghese>



## Color Spaces - YUV (YCbCr)



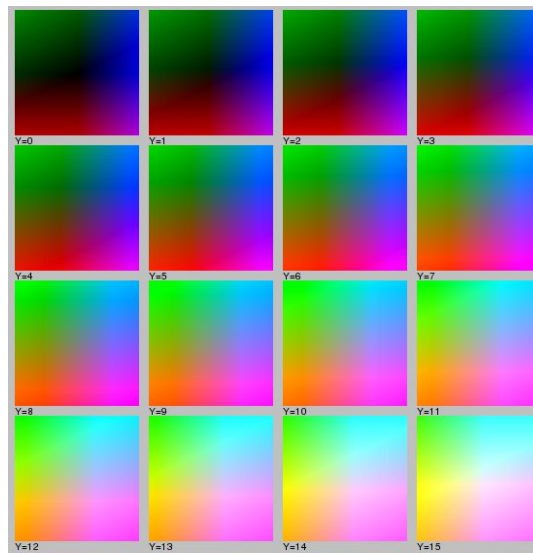
A.A. 2008-2009

25/80

<http://homes.dsi.unimi.it/~borghese>



## Color Spaces - YCbCr



A.A. 2008-2009

26/80

<http://homes.dsi.unimi.it/~borghese>



## Image RGB



A.A. 2008-2009

28/80

<http://homes.dsi.unimi.it/~borghese>



## Color Spaces - Discussion



- RGB
  - ◆ Handled by most capture cards
  - ◆ Used by computer monitors
  - ◆ Not easily separable channels
  
- YCbCr (YUV)
  - ◆ Handled by most capture cards
  - ◆ Used by TVs and JPEG images
  - ◆ Easily workable color space

A.A. 2008-2009

30/80

<http://homes.dsi.unimi.it/~borghese>



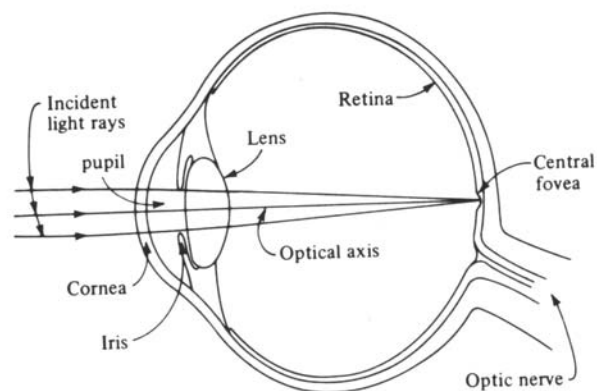
## Sommario



- La visione
- Le immagini digitali
- **Il modello geometrico di una camera**
- Segmentazione real-time.



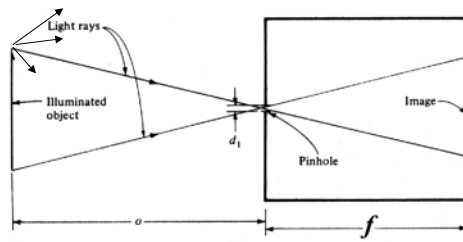
## L'occhio umano



Its behavior is very similar to that of a camera



## La camera come strumento di ripresa



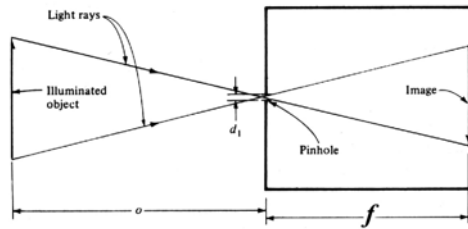
A.A. 2008-2009

34/80

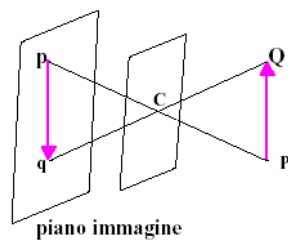
<http://homes.dsi.unimi.it/~borghese>



## La pin-hole camera



**Proiezione prospettica:**  
tutti i raggi di proiezione  
passano per un unico punto,  
detto **centro di proiezione**.



Pinhole camera

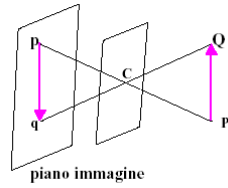
A.A. 2008-2009

35/80

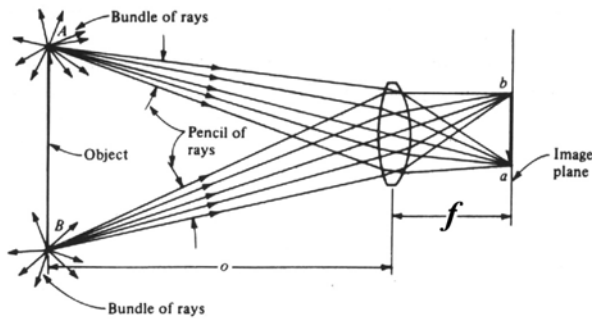
<http://homes.dsi.unimi.it/~borghese>



# La lente



Pinhole camera



Lente convergente

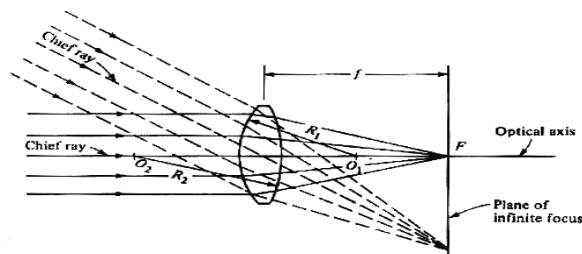
A.A. 2008-2009

36/80

<http://homes.dsi.unimi.it/~borghese>



# Geometria dell'ottica



Oggetti all'infinito

- **Distanza focale:** distanza del piano immagine quando un oggetto si trova all'infinito.
- **Asse ottico:** raggio che non viene deviato dalla lente.
- **Intersezione dell'asse ottico con il piano immagine dà il punto principale (F).**

A.A. 2008-2009

37/80

<http://homes.dsi.unimi.it/~borghese>

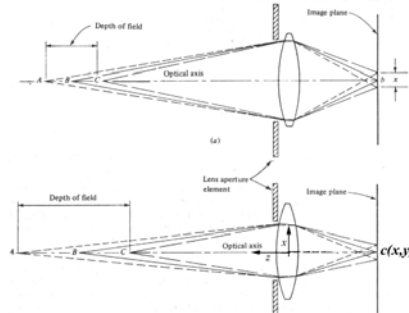




# Messa a fuoco



Problema della messa a fuoco



**Parametri di camera (o intrinseci):**

- Punto principale  $c(x,y)$  + lunghezza focale,  $f$  (3 parametri).
- Occorre conoscere anche il fattore di forma dei pixel nel caso di immagini digitali (è una costante, non un parametro).
- (Distorsioni).

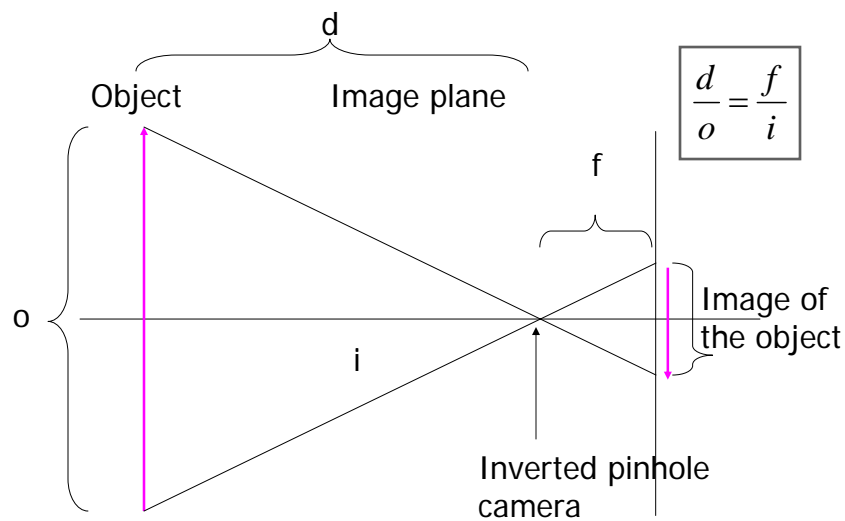
A.A. 2008-2009

38/80

<http://homes.dsi.unimi.it/~borghese>




# Calculating Distance




A.A. 2008-2009

39/80

<http://homes.dsi.unimi.it/~borghese>

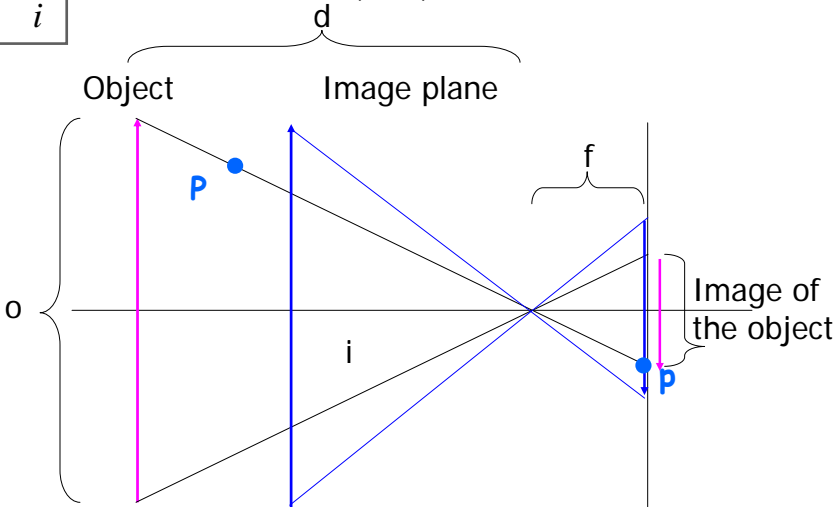


## Observations



The image becomes bigger as the distance from F decreases.  
Any point P on the same projecting line, is projected into the same point p


$$\frac{d}{o} = \frac{f}{i}$$




A.A. 2008-2009

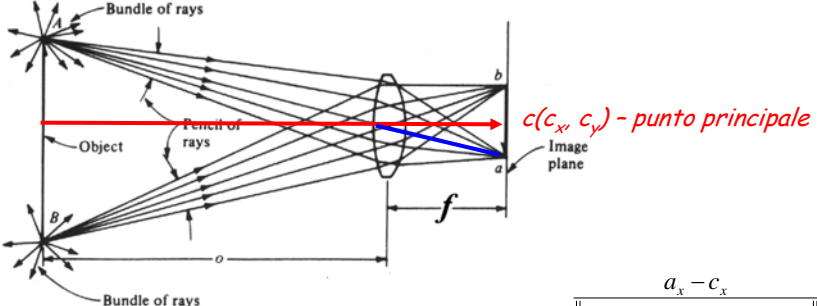
40/80

<http://homes.dsi.unimi.it/~borghese>



## Retroproiezione





$c(c_x, c_y)$  - punto principale

Il punto A si trova lungo la retta che ha coseni direttori, misurati in un sistema di riferimento locale:

$$\frac{a_x - c_x}{\|f + (a_x - c_x) + (a_y - c_y)\|}$$

$$\frac{a_y - c_y}{\|f + (a_x - c_x) + (a_y - c_y)\|}$$

$$\frac{f}{\|f + (a_x - c_x) + (a_y - c_y)\|}$$

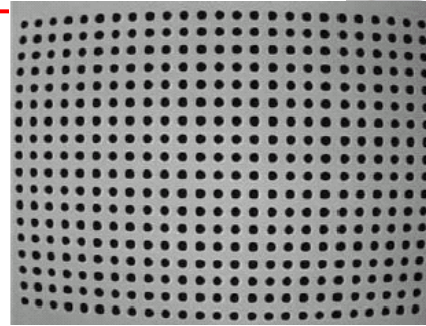
A.A. 2008-2009

41/80

<http://homes.dsi.unimi.it/~borghese>



## Esempi di Distorsioni



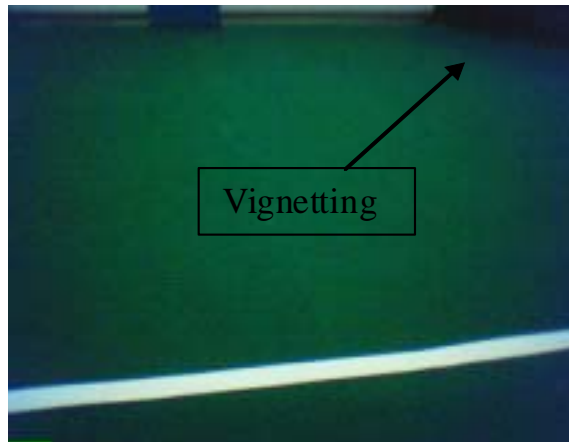
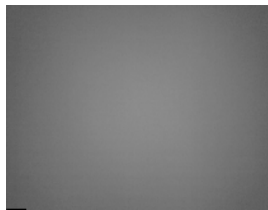
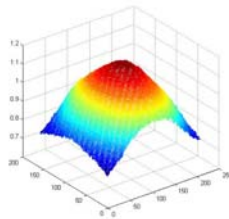
Ottime per effetti speciali, un po' meno per delle misure.....  
*Le camere non sono metriche.*

A.A. 2008-2009

42/80



## Vignetting



Compensazione della risposta del sensore

A.A. 2008-2009

43/80

<http://homes.dsi.unimi.it/~borghese>



## Sommario



- La visione
- Le immagini digitali
- Il modello geometrico di una camera
- **Segmentazione real-time.**



## Fast color segmentation



- Low level vision is responsible for summarizing *relevant-to-task* image features
  - ◆ Color is the main feature that is relevant to identifying the objects needed for the task
  - ◆ Important to reduce the total image information
- **Color segmentation algorithm**
  - ◆ Segment image into *symbolic colors*
  - ◆ Run *length encode* image
  - ◆ Find *connected components*
  - ◆ Join nearby components into *regions*



## Color Segmentation



- Goal: semantically label each pixel as belonging to a particular type of object
- Map the domain of raw camera pixels into the range of symbolic colors  $L$ 
  - ◆  $C$  may include for instance ball, carpet, 2 goal colors, 1 additional marker color, 2 robot colors, walls/lines and unknown

$$F : y, u, v \rightarrow c \in L$$

- Reduces the amount of information per pixel roughly by 1.8M
  - ◆ Instead of a space of  $2^{24}$  (R,G,B each on 8 pixels) =  $256^3$  values, we only have 9 values!

NB No semantic is given. Low level vision.



## Before Segmentation





## Ideal Segmentation



A.A. 2008-2009

48/80

<http://homes.dsi.unimi.it/~borghese>



## Result of Segmentation



Segmentazione HW in 8 classi. Viene definita per ogni classe:

$$\begin{array}{ll} Y_{\min} & Y_{\max} \\ Cb_{\min} & Cb_{\max} \\ Cr_{\min} & Cr_{\max} \end{array}$$

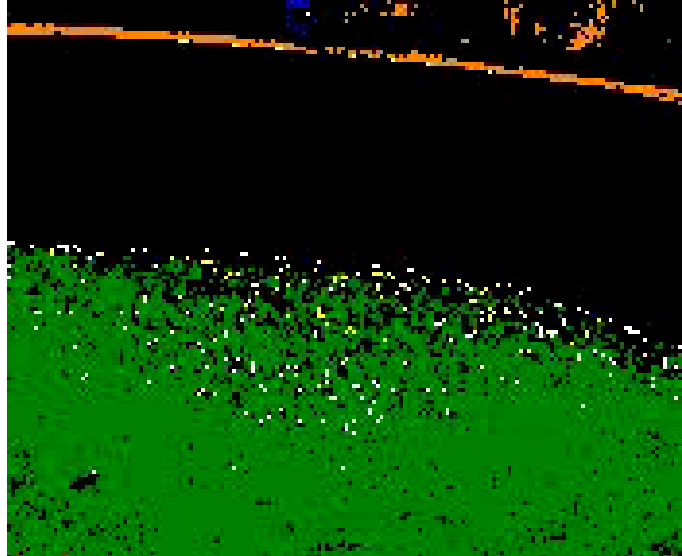
A.A. 2008-2009

49/80

<http://homes.dsi.unimi.it/~borghese>



## Potential Problems with Color Segmentation



A.A. 2008-2009

50/80

<http://homes.dsi.unimi.it/~borghese>



## Color Segmentation Analysis

- Advantages
  - ◆ Quickly extract relevant information
  - ◆ Provide useful representation for higher-level processing
  - ◆ Differentiate between YCbCr pixels that have *similar* values
- Disadvantages
  - ◆ Cannot segment YCbCr pixels that have *identical* values into different classes
  - ◆ Generate smoothly contoured regions from noisy images
- HW implementation
  - ◆ Multi-thresholding to extract the different classes.

A.A. 2008-2009

51/80

<http://homes.dsi.unimi.it/~borghese>

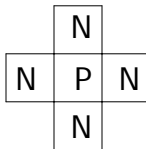


## Turning Pixels into Regions

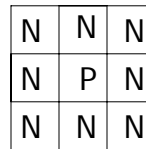
- A disjoint set of labeled pixels is still not enough to properly identify objects
- Pixels must be grouped into spatially-adjacent regions
  - ◆ Regions are grown by considering local neighborhoods around pixels

How to achieve this?

4-connected neighborhood



8-connected neighborhood



A.A. 2008-2009

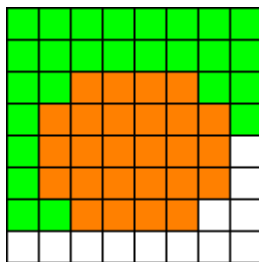
52/80

<http://homes.dsi.unimi.it/~borghese>

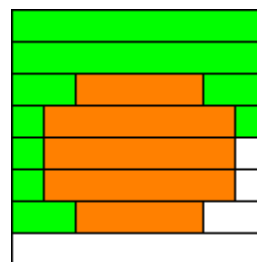


## 1) Run Length Encoding

- Segment each image row into groups of similar pixels called *runs*
  - ◆ Runs store a start and number of pixels belonging to the **same** color (alternatively the end point is memorized)



Original image



RLE image

A.A. 2008-2009

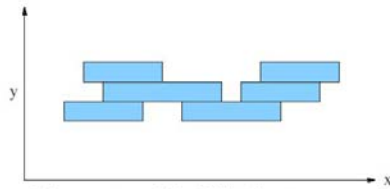
53/80

<http://homes.dsi.unimi.it/~borghese>

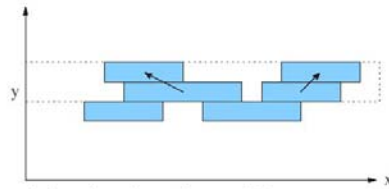




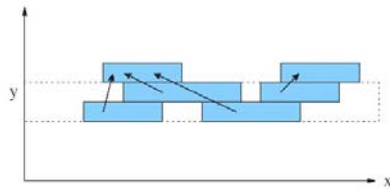
## 2) Merging Regions



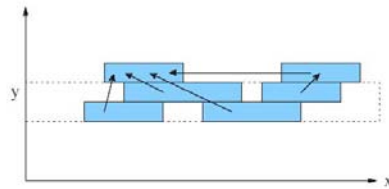
1: Runs start as a fully disjoint forest



2: Scanning adjacent lines, neighbors are merged



3: New parent assignments are to the furthest parent



4: If overlap is detected, latter parent is updated

Goal: a closed region is formed. Isolated pixels are discarded.

A.A. 2008-2009

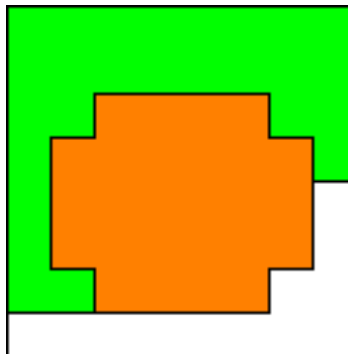
54/80

<http://homes.dsi.unimi.it/~borghese>



## Final Results

- Runs are merged into multi-row regions
- Image is now described as contiguous regions instead of just pixels
- Objects are convexes



A.A. 2008-2009

55/80

<http://homes.dsi.unimi.it/~borghese>



## Data Extracted from Regions



- Features extracted from regions
  - ◆ *Centroid*
    - ☞ Mean location
  - ◆ *Bounding box*
    - ☞ Max and min (x,y) values
  - ◆ *Area*
    - ☞ Number of pixels in box
  - ◆ *Average color*
    - ☞ Mean color of region pixels
  
- Regions are stored by color class and sorted by largest area
- These features let us write concise and fast object detectors

Features are a function of later processing required.



## High level vision - Object Detection Process



- Produces a set of candidate objects that might be this object from lists of regions
  - ◆ Given 'n' orange blobs, is one of them the ball?
  
- Compares each candidate object to a set of **models** that predict what the object would look like when seen by a camera
  - ◆ **Models** encapsulate all assumptions
  - ◆ Also called filtering
  
- Selects best match to report to behaviors
  - ◆ Position and quality of match are also reported



## Filtering Overview



- Each filtering **model** produces a number in  $[0.0, 1.0]$  representing the certainty of a match
  - ◆ Some filters can be binary and will return either 0.0 or 1.0
- Certainty levels are multiplied together to produce an overall match
  - ◆ Real-valued range allows for areas of uncertainty
  - ◆ Keeps one bad filter result from ruining the object
  - ◆ Multiple bad observations will still cause the object to be thrown out

Alternatively fuzzy inference can be used.



## Object Detection Algorithm



- 1) Produce a set of candidate objects that might be the object from the list of regions produced by low-level vision.
- 2) Compare each candidate object to a set of models that predict features that the object should have when seen through a camera.
- 3) Select the best match and compute its distance to the robot based on the **robot's camera model and kinematics**.
- 4) Return this best match along with the quality of the estimate to the behaviors.



# Riconoscimento palla - CMU - I



- Minimum size
  - Makes sure the ball has a bounding box at least 3 pixels tall and wide and 7 pixels total area
- Square bounding box
  - Makes sure the bounding box is roughly square
  - Uses an unnormalized Gaussian as the output
  - Output is as follows:

$$d = \frac{w-h}{w+h} \quad \text{Bounding box (height x width)}$$

$$o = e^{-\left(\frac{d}{C}\right)^2 / 2} \quad \begin{matrix} C=0.2 \text{ if on edge of image} \\ 0.6 \text{ otherwise} \end{matrix}$$

o = 1 iff w = h  
o < 1 if w > h or w < h

15-491 CMRoboBits



# Riconoscimento palla - CMU - II

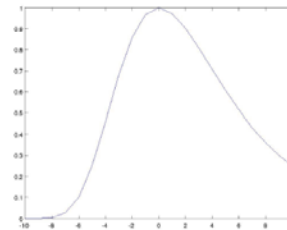


Filter

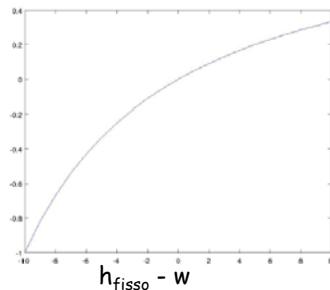
$$d = \frac{w-h}{w+h}$$

$$o = e^{-\left(\frac{d}{C}\right)^2 / 2}$$

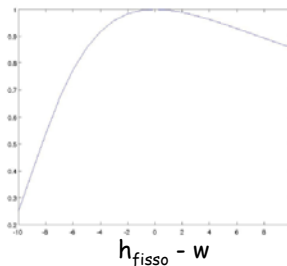
Plot: o  
C=0.2



Plot: d  
H=10  
W=[0-20]



Plot: o  
C=0.6



**Calculating Distance**

The diagram illustrates an inverted pinhole camera. An object of height  $o$  is located at a distance  $d$  from the image plane. The image plane is at a distance  $i$  from the camera's pinhole. The focal length of the camera is  $f$ . The image of the object is formed on the image plane. The relationship between these distances is given by the equation:

$$\frac{d}{o} = \frac{f}{i}$$

A.A. 2008-2009 68/80 http://homes.dsi.unimi.it/~borghese

**Calculation of Camera Position**

- Position of camera is calculated based on body position and head position w.r.t body
- Body position is known from walk engine
- Head position relative to body position is found from forward kinematics using joint positions
- Camera position
  - ◆ *camera\_loc* is defined as position of camera relative to egocentric origin
  - ◆ *camera\_dir*, *camera\_up*, and *camera\_down* are unit vectors in egocentric space
  - ↳ Specify camera direction, up and right in the image

The diagram shows a 3D model of a robot's head and neck assembly. The camera position is labeled as *camera position (camera\_loc)*. Unit vectors for camera direction (*camera\_dir*), camera up (*camera\_up*), and camera right (*camera\_right*) are shown. The origin for egocentric coordinates is marked at the base of the neck. A *body angle* is indicated relative to the body height.

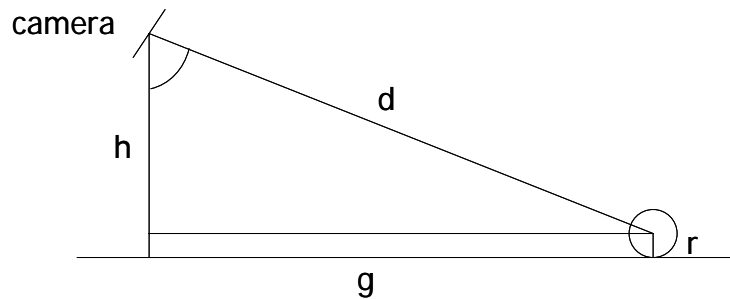
A.A. 2008-2009



## Ball Position Estimation



- Two methods are used for estimating the position of the ball
  - The first calculates the camera angle from the ball model
  - The second uses the robot's encoders to calculate the head angle
- The first is more accurate but relies on the pixel size of the ball
  - This method is chosen if the ball is NOT on the edge of the image
  - Partial occlusions will make this estimate worse



A.A. 2008-2009

70/80

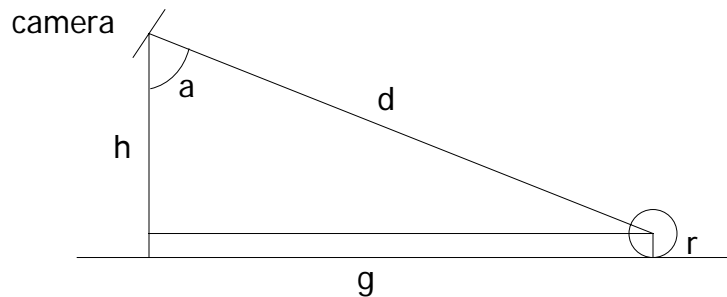
<http://homes.dsi.unimi.it/~borghese>



## Ball Position Estimation - I It works also for the robotic arm



- Ball position estimation problem is overconstrained.
  - $g$  is the unknown
  - Works also when the ball is at the edge of the image
  - Critical estimate of the "body schema" (height  $h$  and angle  $a$ )
  - $g = h \tan(a)$ , does not need  $d$



A.A. 2008-2009

71/80

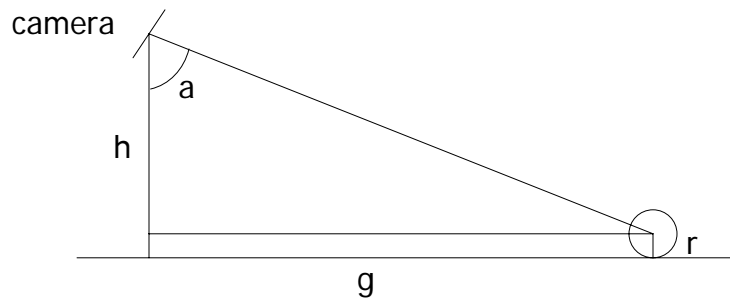
<http://homes.dsi.unimi.it/~borghese>



## Ball - Position Estimation - I



- This method works all of the time
  - ◆ Camera angle computed from kinematics (but it may not be accurate enough!)
  - ◆ Used when ball is on edge of image



A.A. 2008-2009

72/80

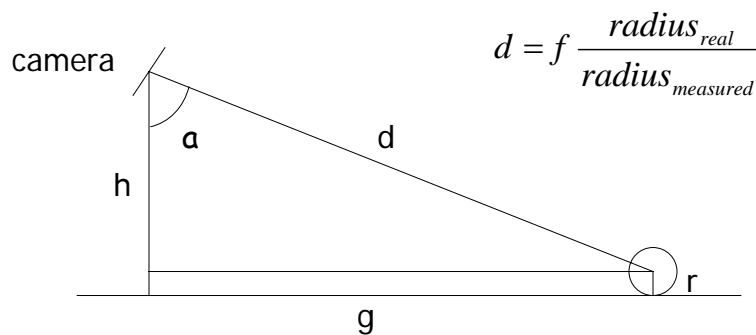
<http://homes.dsi.unimi.it/~borghese>



## Ball - Position Estimation - II



- This method is more accurate
  - ◆ Requires accurate pixel count
  - ◆ Used only when ball is near center of image
  - ◆ More reliable when many pixels are available (close range)
  - ◆  $g = \sqrt{d^2 - h^2}$
  - ◆  $a = \arccos(h/d)$



$$d = f \frac{\text{radius}_{\text{real}}}{\text{radius}_{\text{measured}}}$$

A.A. 2008-2009

73/80

<http://homes.dsi.unimi.it/~borghese>



## End Result - Accurate Ball Position



A.A. 2008-2009

76/80

<http://homes.dsi.unimi.it/~borghese>



## AIBO's vision::OVirtualRobotComm



- Servizio utilizzato per accedere alla telecamera:
  - ◆ `OvirtualRobotComm.FbkImageSensor.OFbkImageVe`  
`ctorData.S`

A.A. 2008-2009

77/80

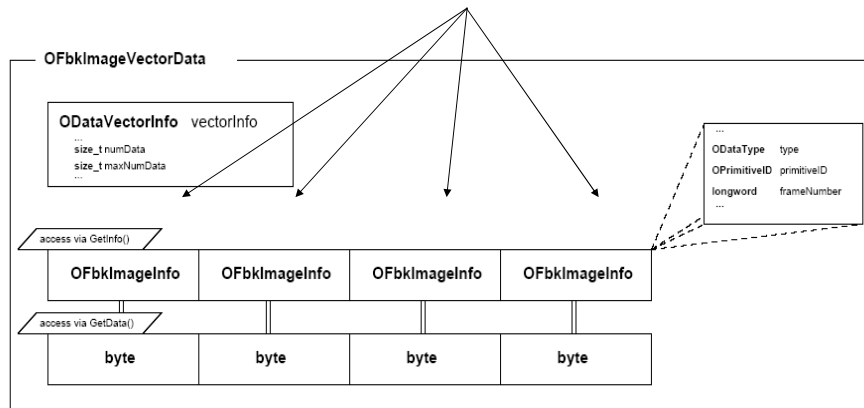
<http://homes.dsi.unimi.it/~borghese>





## Struttura dati telecamera

4 elementi: uno per layer, High, Medium, Low resolution and thresholded.



A.A. 2008-2009

78/80

<http://homes.dsi.unimi.it/~borghese>



## OFbkImage

- `OFbkImage`: oggetto per la gestione delle immagine della telecamera:
  - ◆ Costruttore di `OFbkImage`: `OFbkImage img( imageVector->GetInfo(ofbkimageLAYER_M), imageVector->GetData(ofbkimageLAYER_M), ofbkimageBAND_Y);`
    - ☞ Oggetto per la gestione della banda Y del layer M (medium resolution).
  - ◆ Alcuni metodi utili di `OFbkImage`: `IsValid()`, `Pointer()`, `Width()`, `Height()`, `Pixel(x,y)`...

A.A. 2008-2009

79/80

<http://homes.dsi.unimi.it/~borghese>



## Sommario



- La visione
- Le immagini digitali
- Il modello geometrico di una camera
- Segmentazione real-time.