

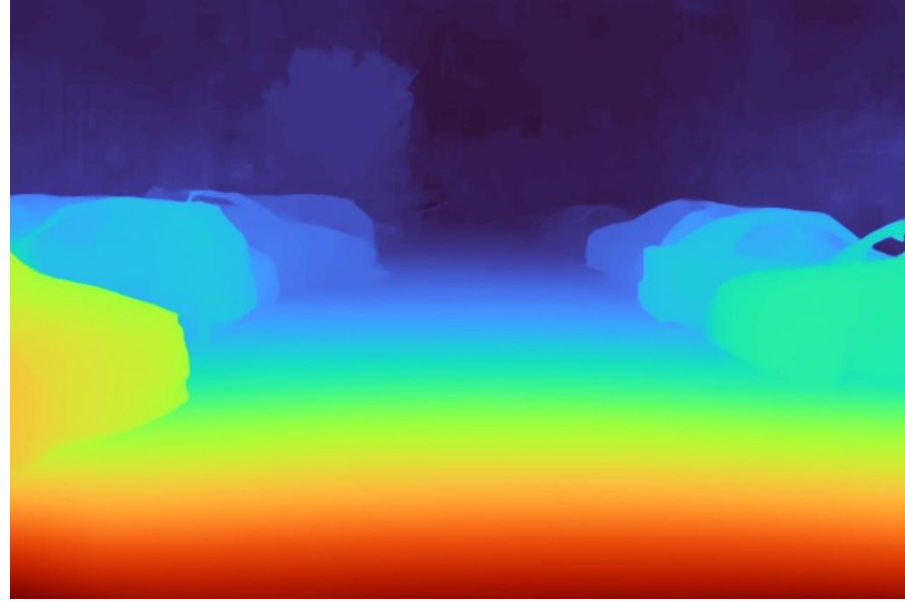
# Semillero de Investigación “Hands - on” Computer Vision



# SESIÓN 5: PASSIVE DEPTH

# Contenidos

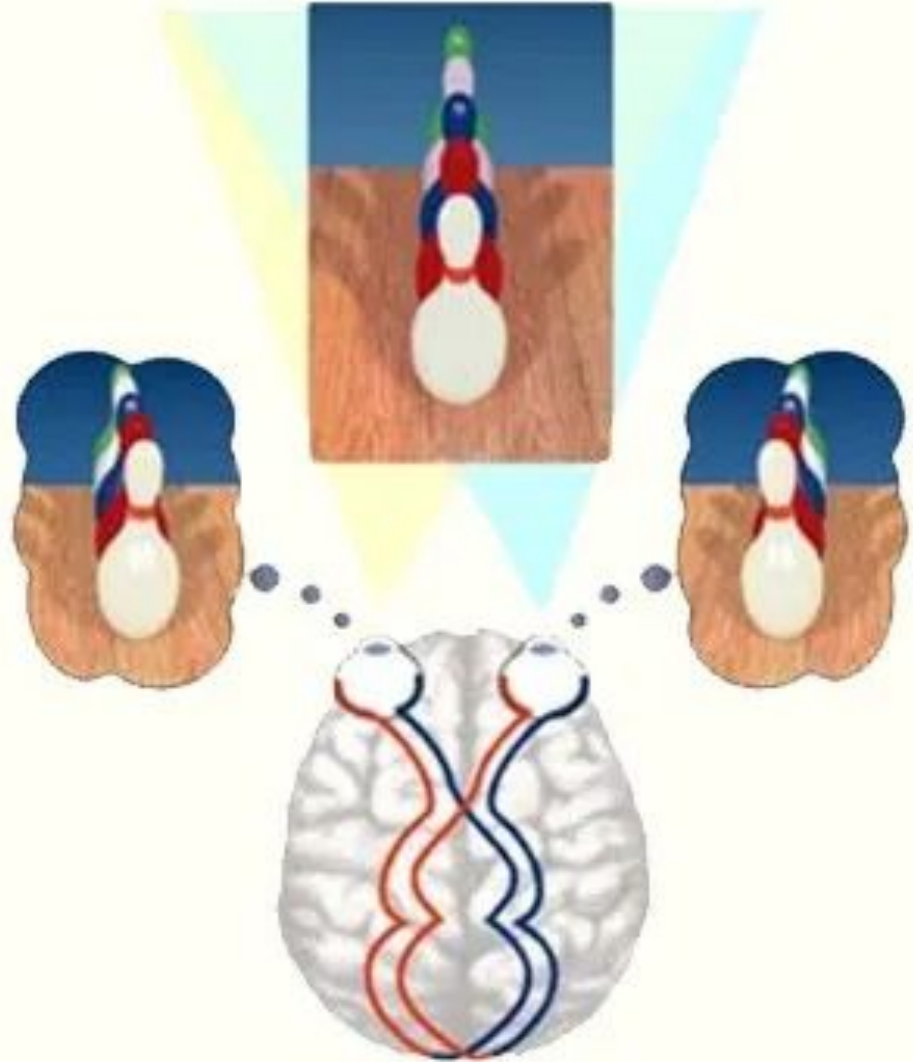
1. Percepción de profundidad
2. Técnicas de estimación de profundidad
  - a. Métodos pasivos
  - b. Métodos activos
3. Visión stereo
  - a. Pipeline
  - b. Disparidad y profundidad
  - c. Rectificación de cámaras
  - d. Triangulación y estimación 3D
4. Hands-on Stereo Vision



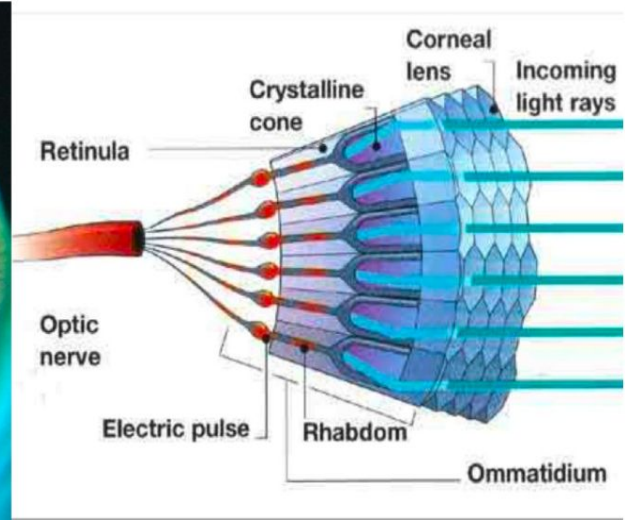
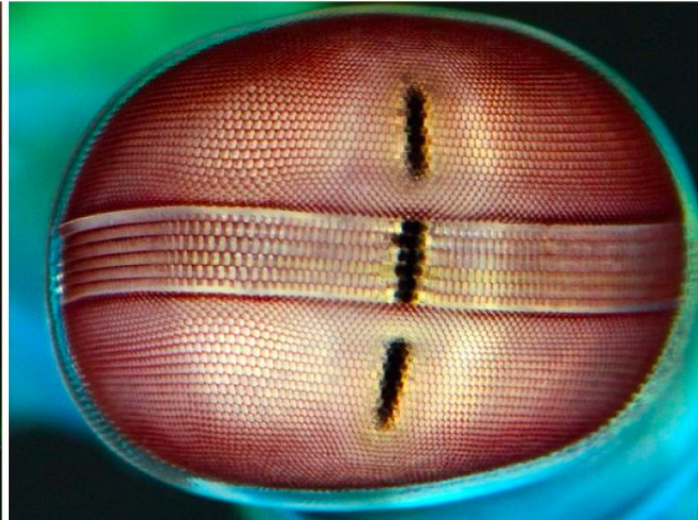
# 1. Percepción de profundidad



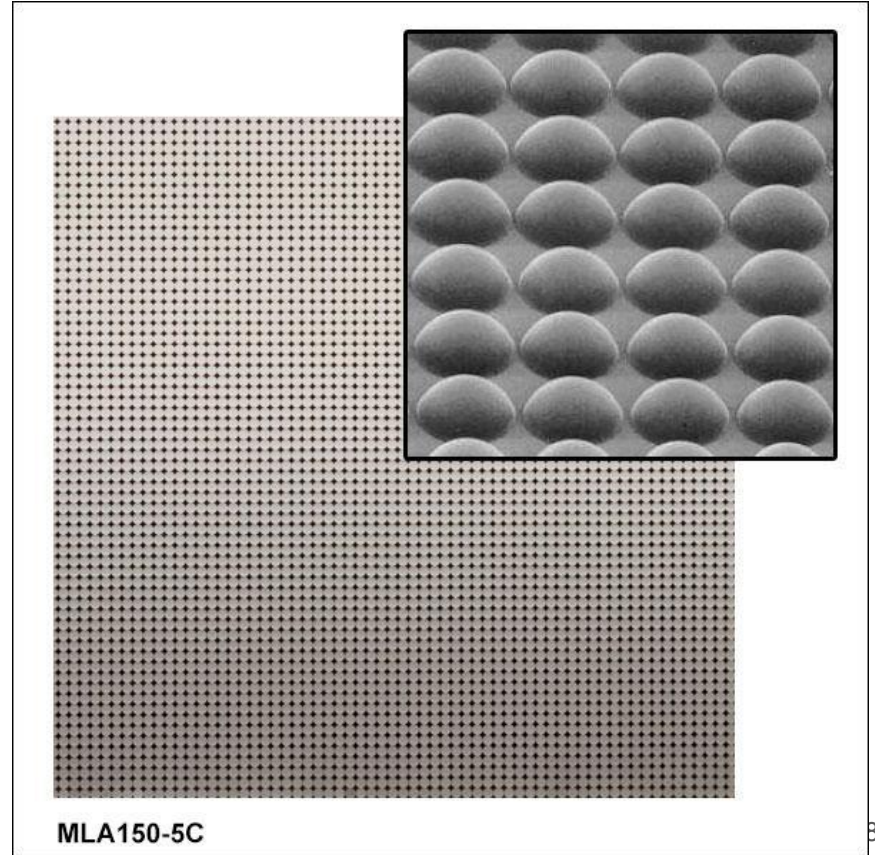
# Percepción de profundidad



# Más allá de Visión Binocular

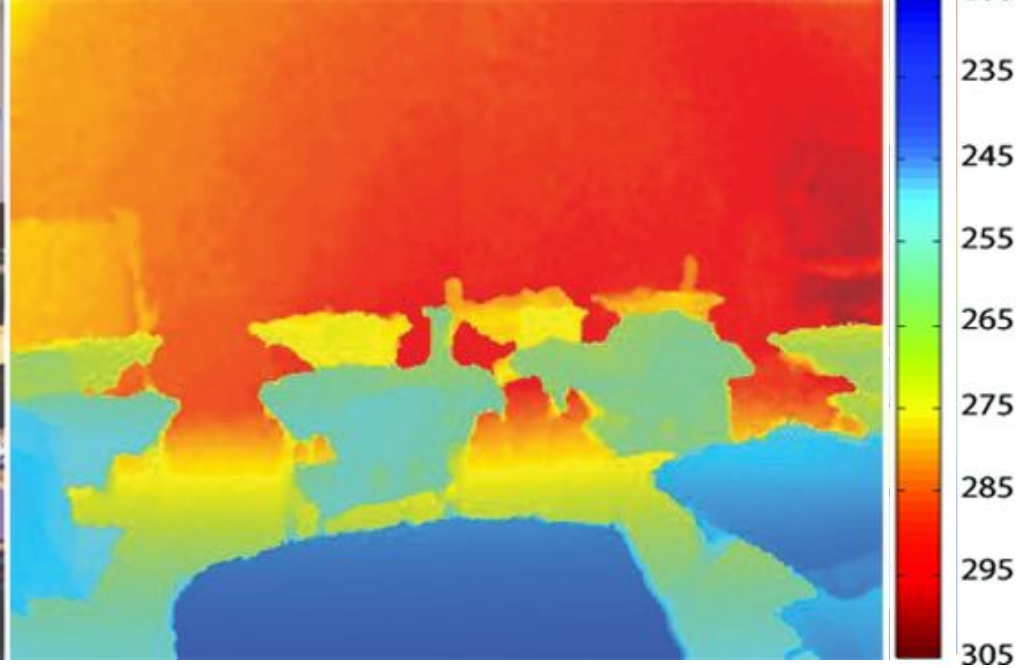


# Múltiples “pares” de cámaras



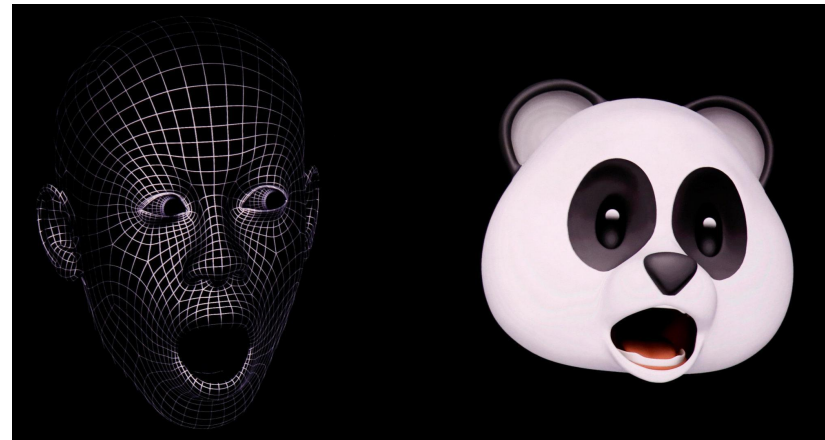
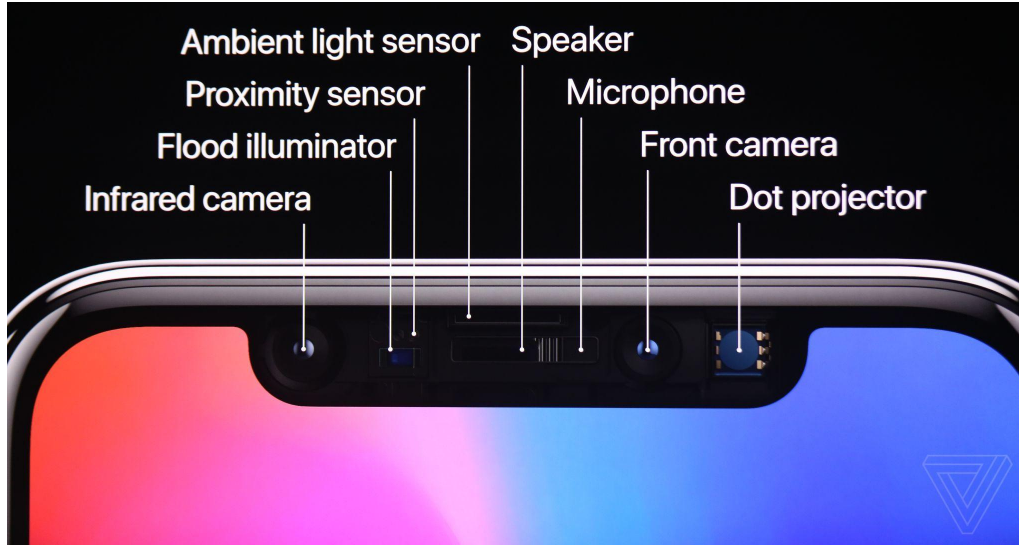


# Imágenes de profundidad (x,y,z)

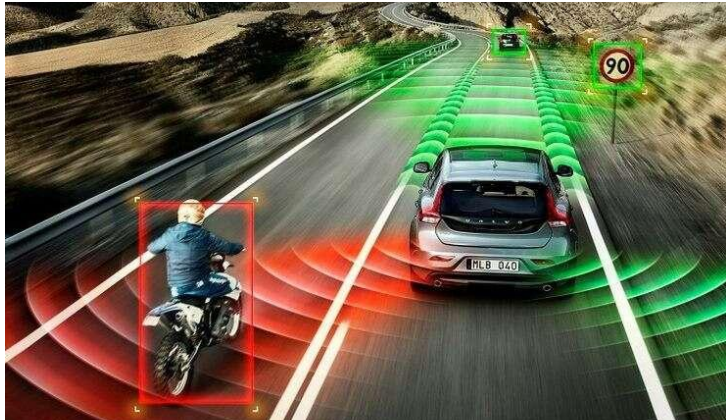
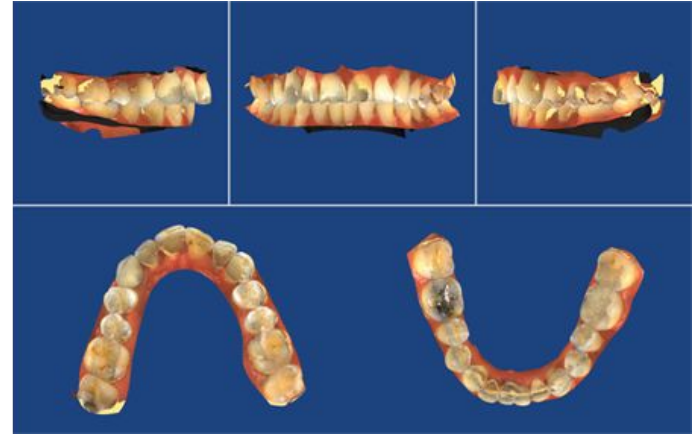
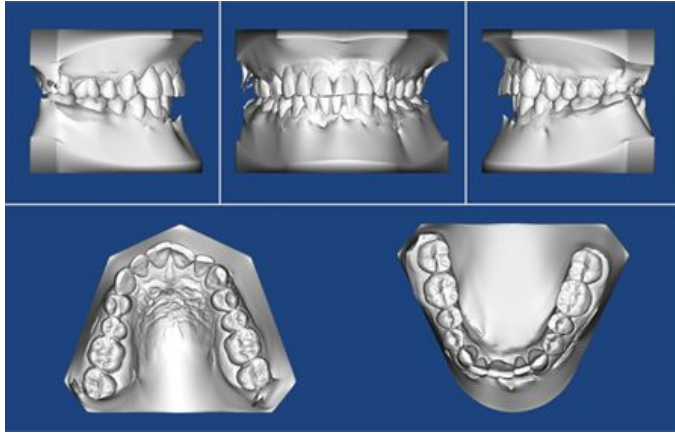


[cm]

# Aplicaciones



# Aplicaciones



**2D**



**3D**



**4D**



**HD**



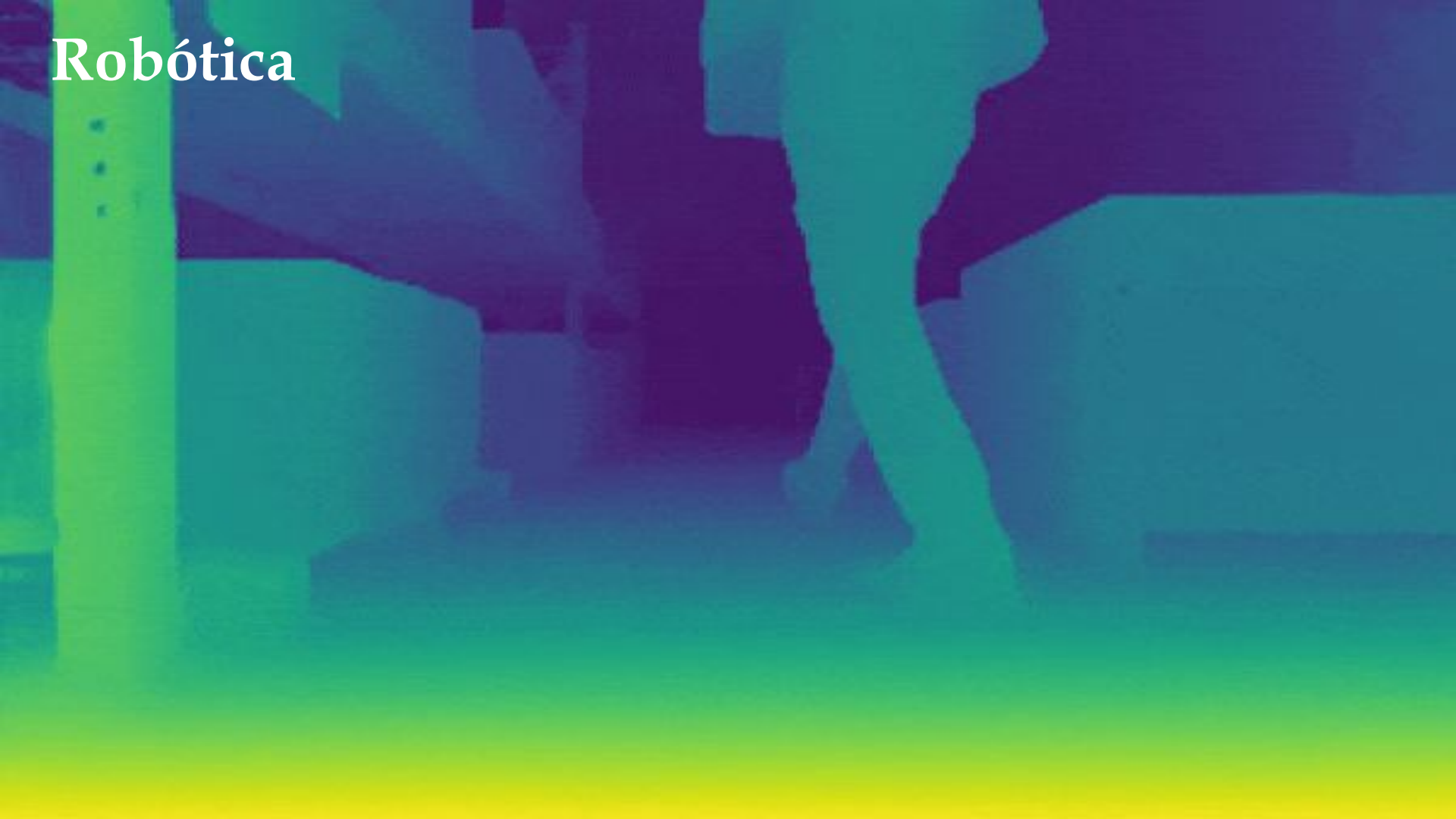


Input



Our depth predictions\*

# Robótica





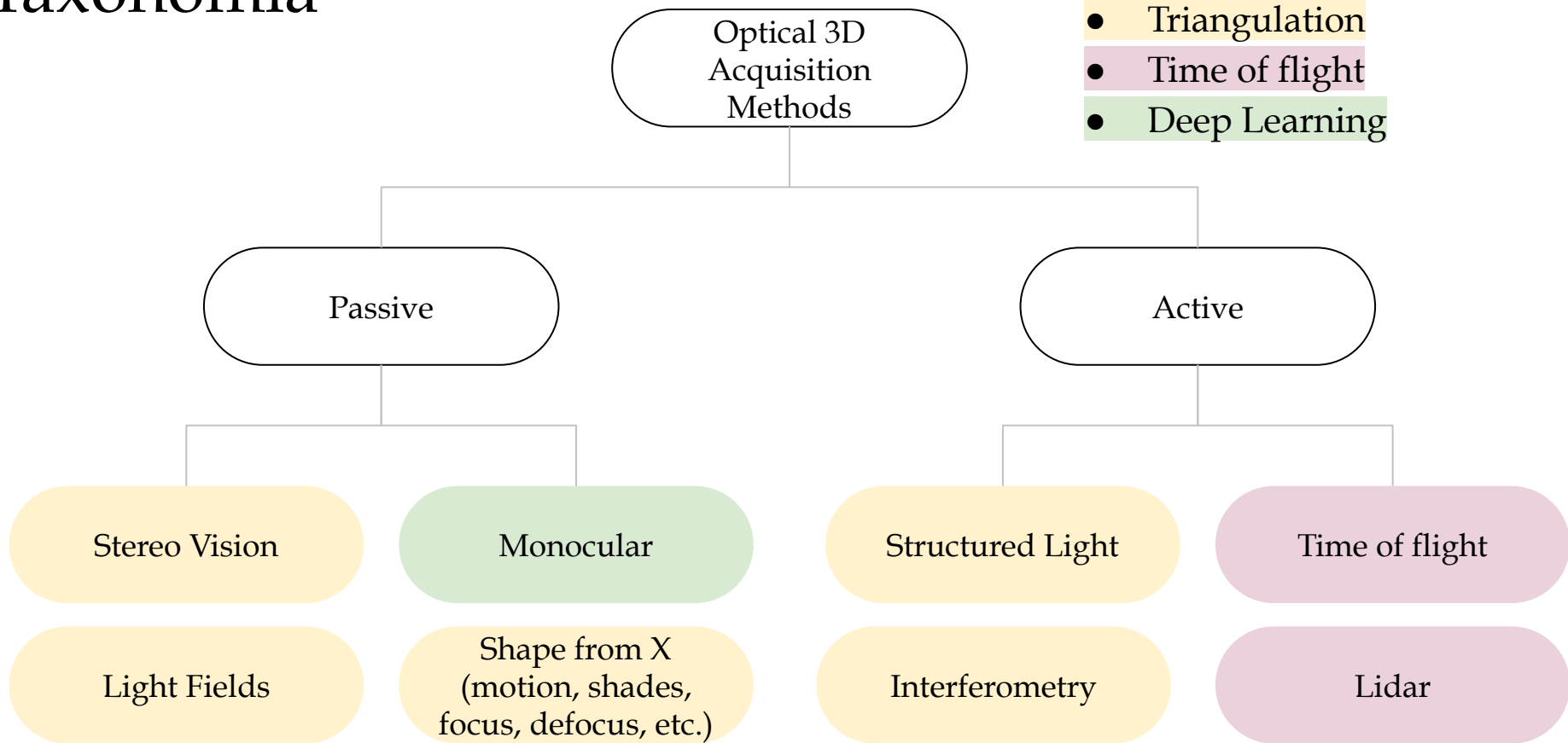
**Vehículos autónomos**

## 2. Técnicas de estimación de profundidad

# Taxonomía

Depth estimation techniques:

- Triangulation
- Time of flight
- Deep Learning

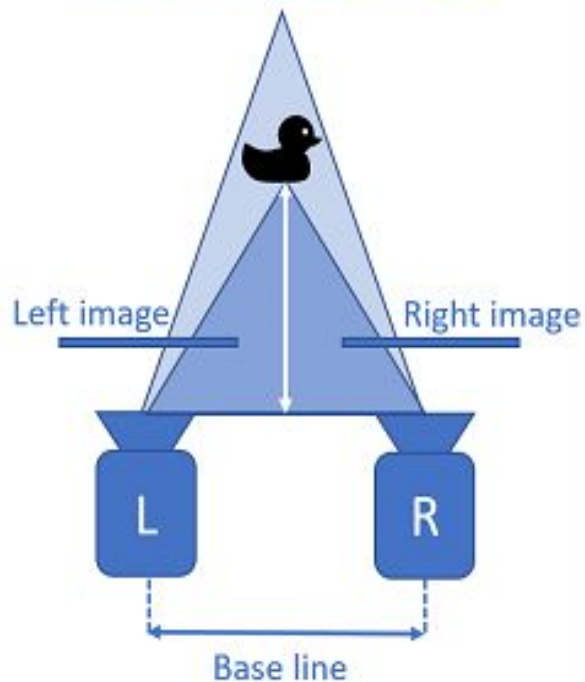


**Pasivo:** Estimación de la profundidad no requiere la emisión activa de señales o radiación hacia la escena

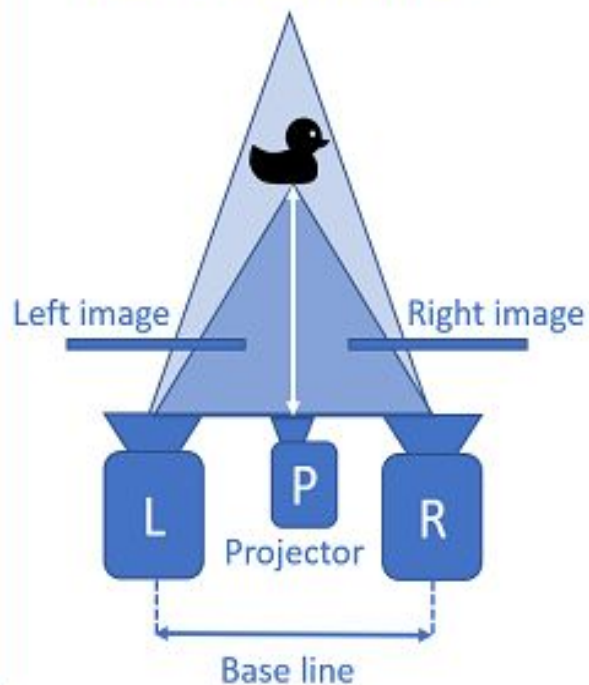


# Basadas en Triangulación

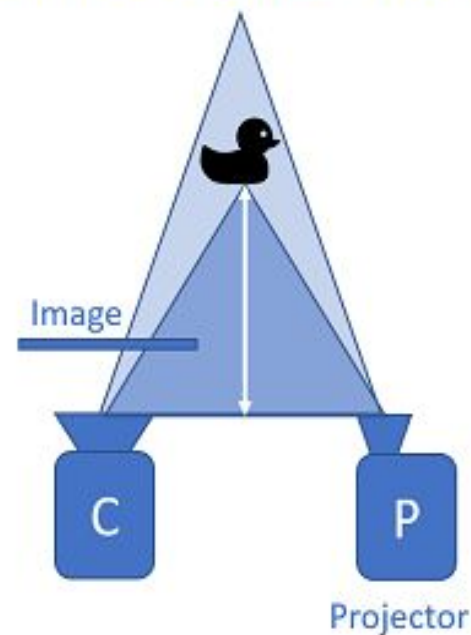
## PASSIVE STEREO



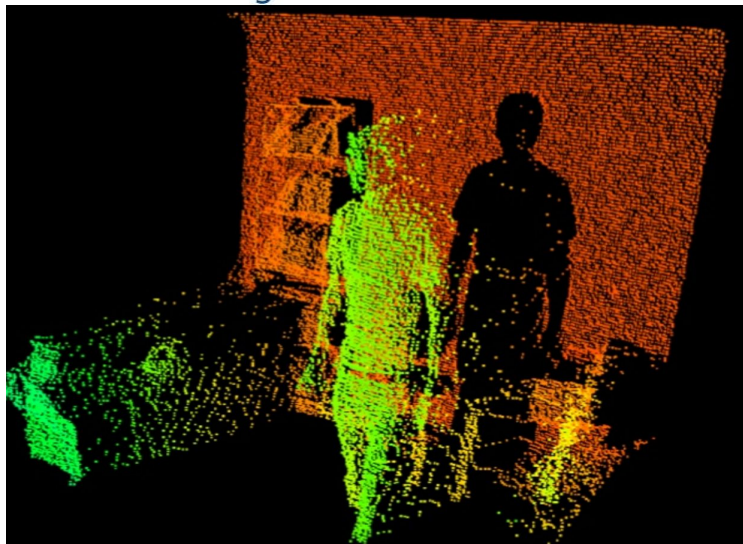
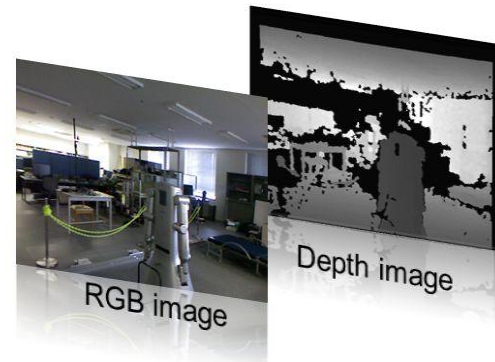
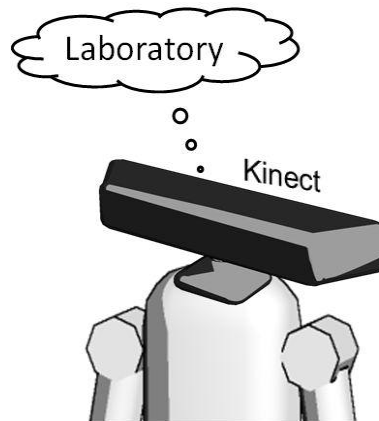
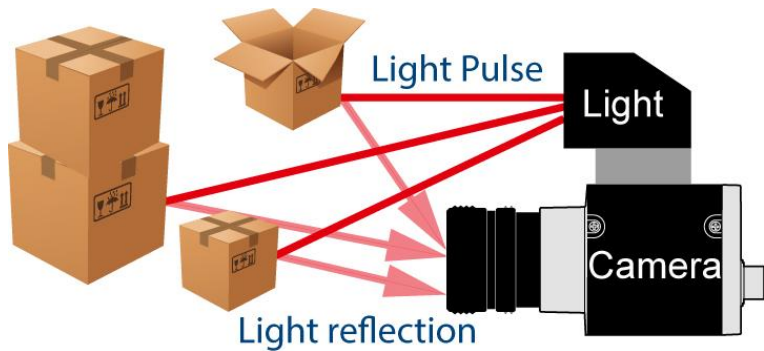
## ACTIVE STEREO



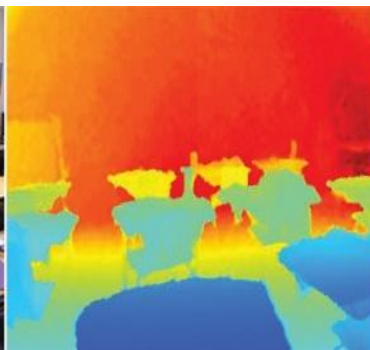
## STRUCTURED LIGHT



# Basadas en Tiempo de Vuelo (Next HoCV session)



RGB



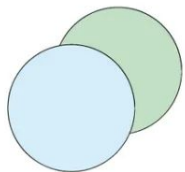
Depth map



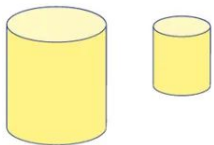
Segmentation

# Basadas en Deep Learning

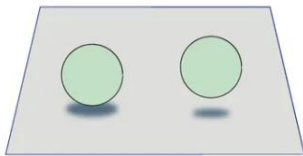
## Monocular Visual Cues



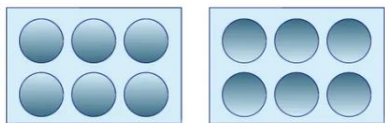
Occlusion



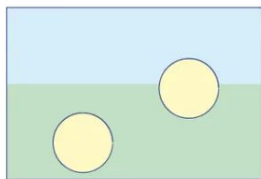
Relative size



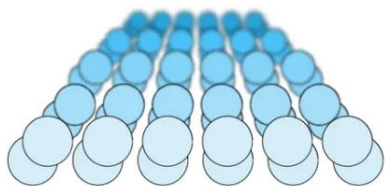
Cast Shadows



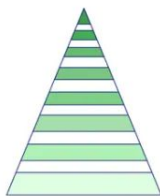
Shading



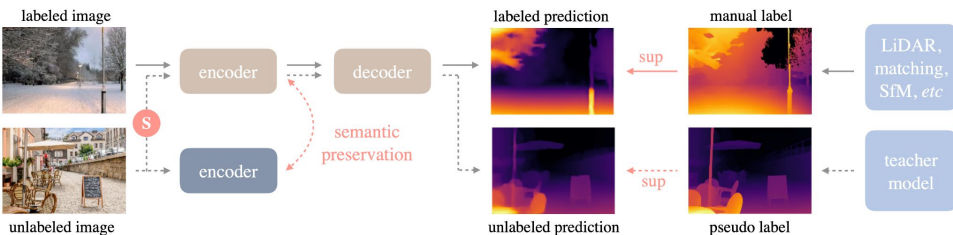
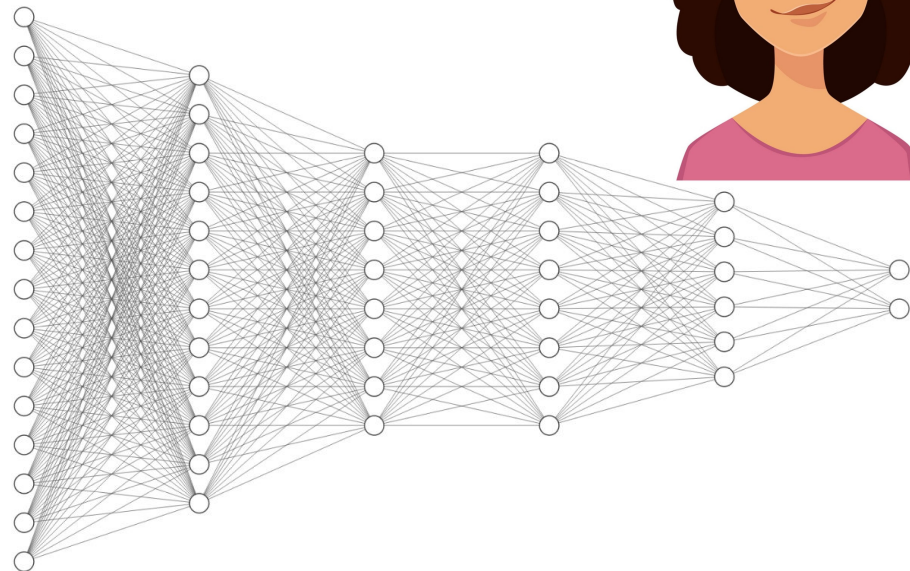
Distance to horizon



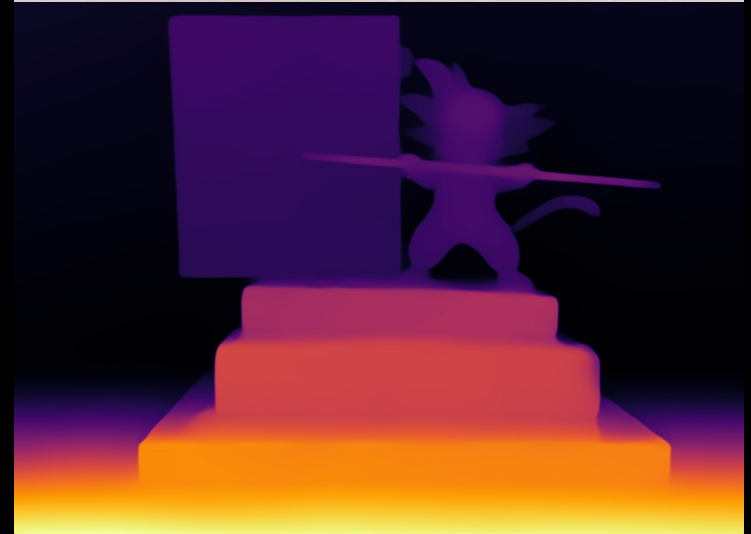
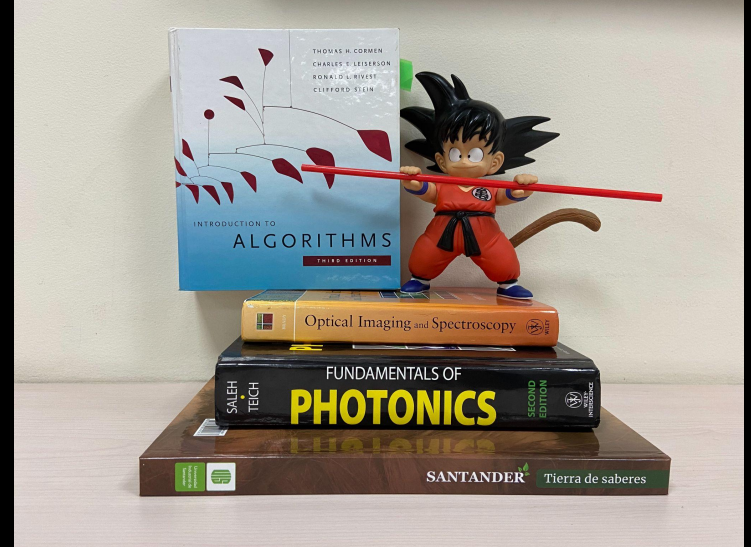
Texture gradient



Linear perspective



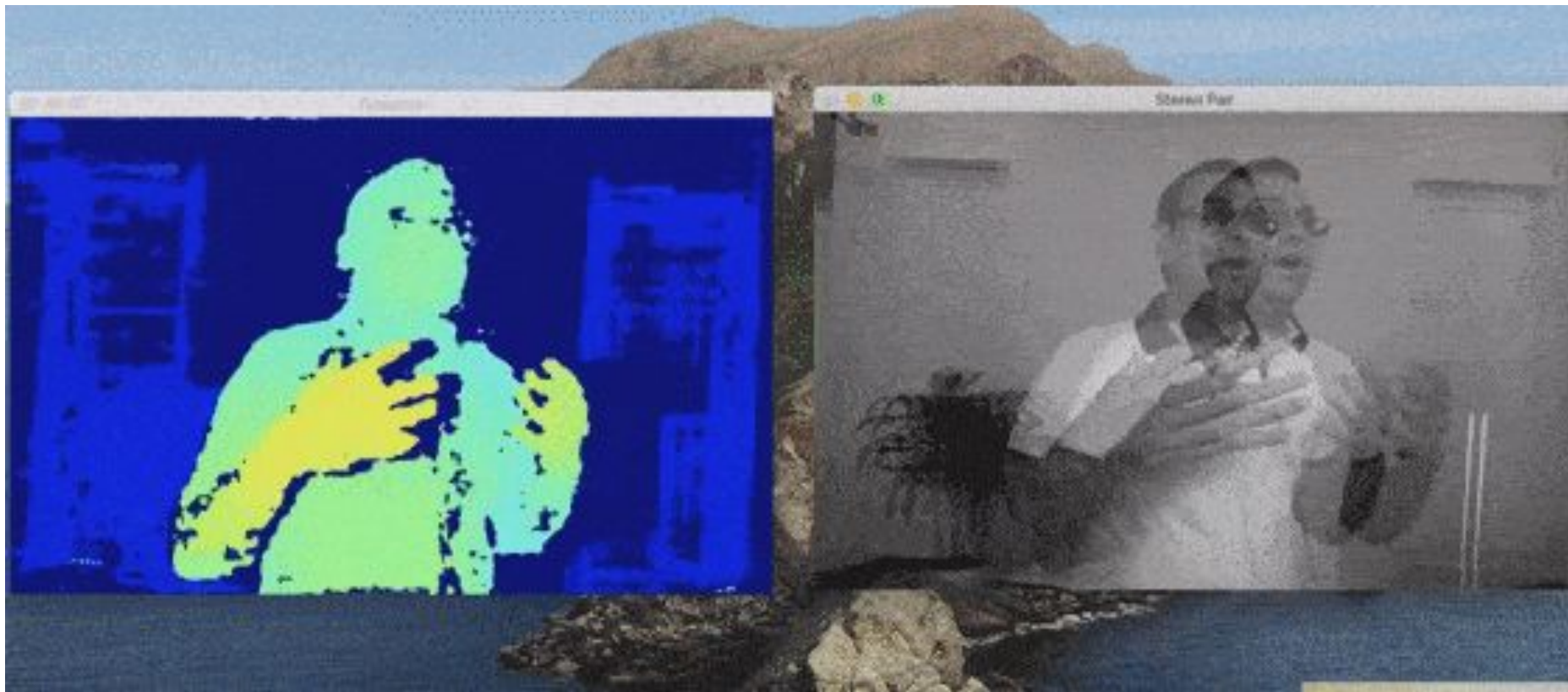
# DEMO (Hugging Face)



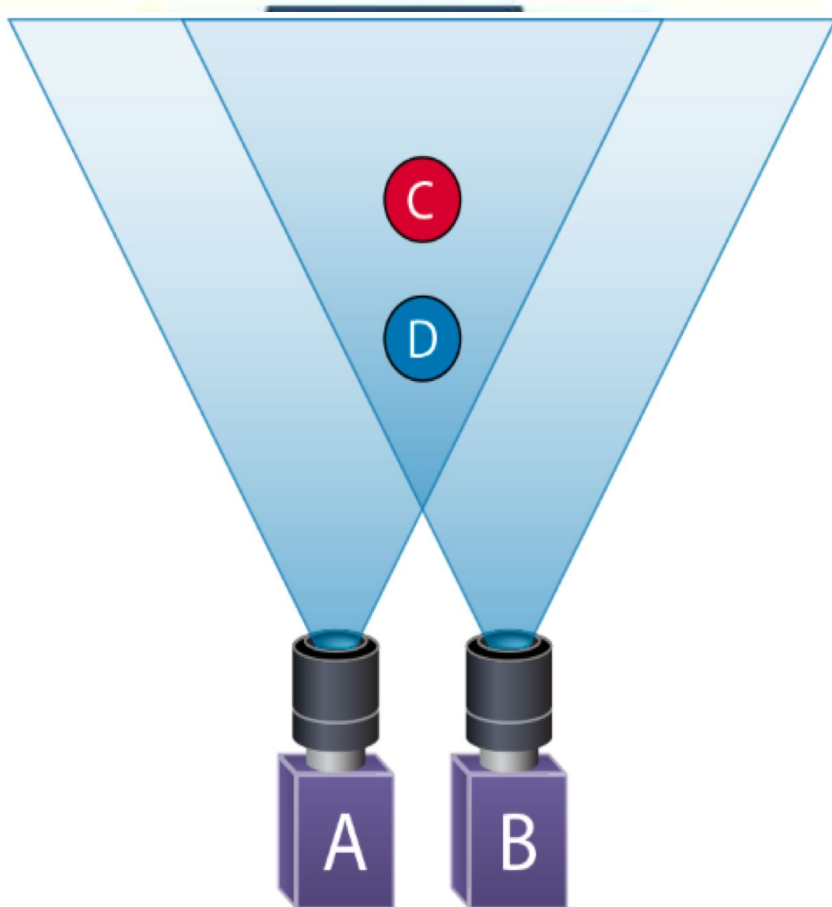
### 3. Visión Stereo (pasiva)



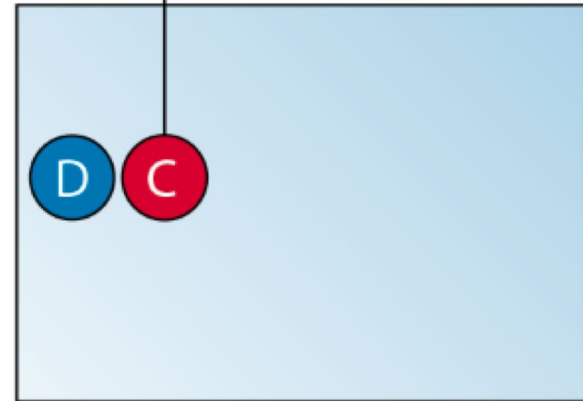
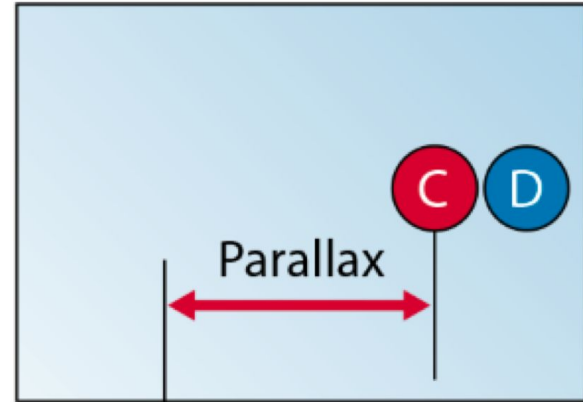
# Visión Stereo (Estimar profundidad a partir de 2 imágenes)



# Parallax (Disparidad)

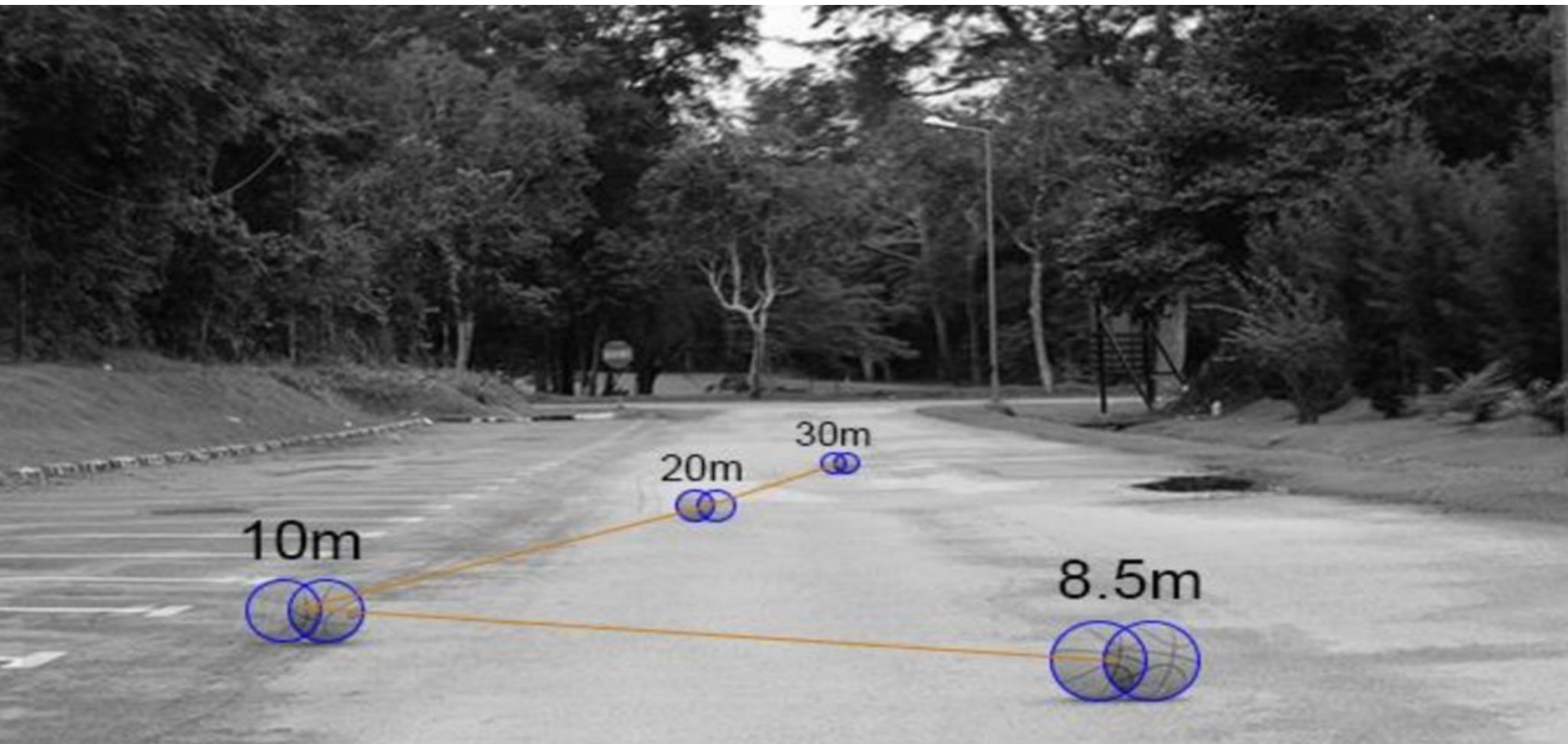


What camera A sees



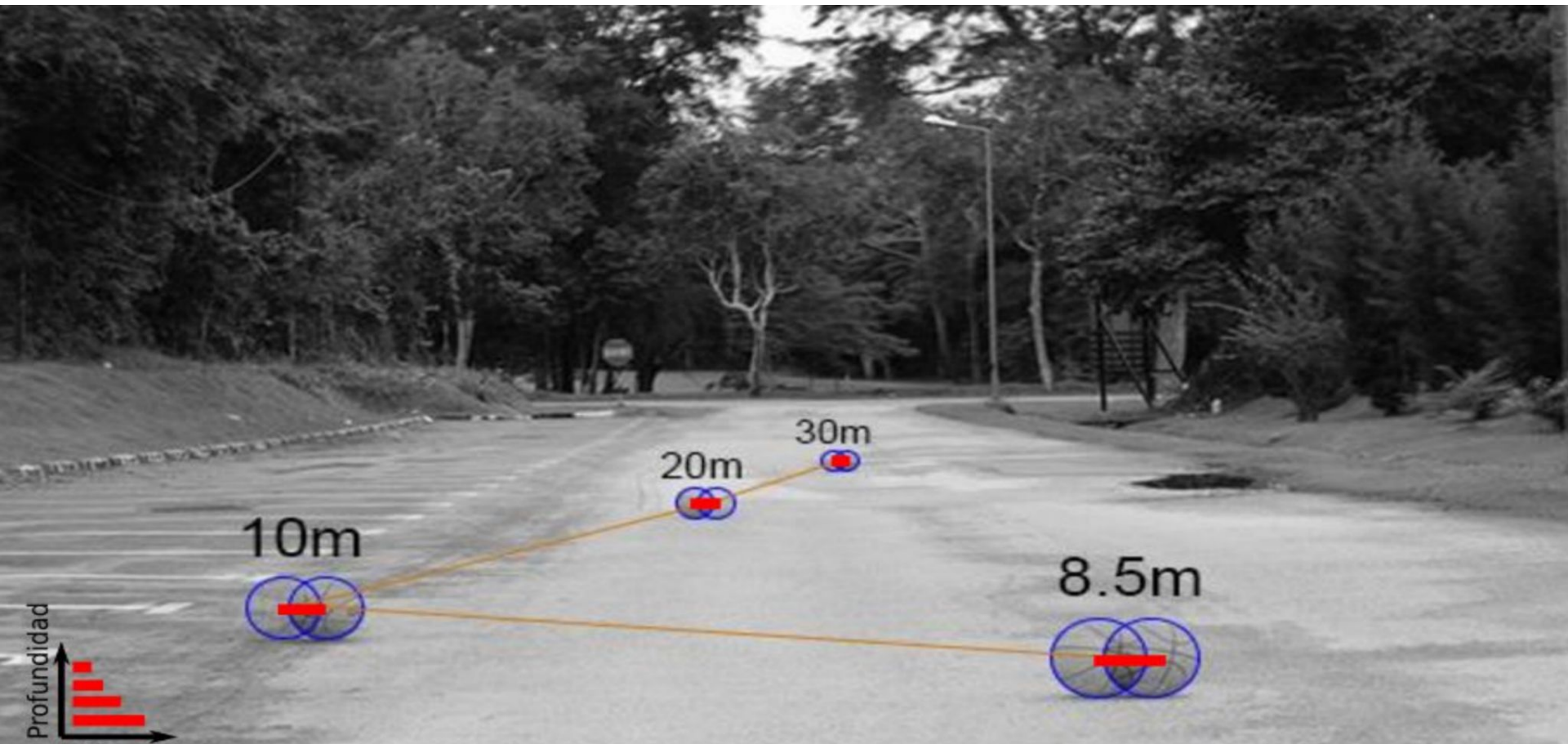
What camera B sees

La disparidad **disminuye** con la distancia

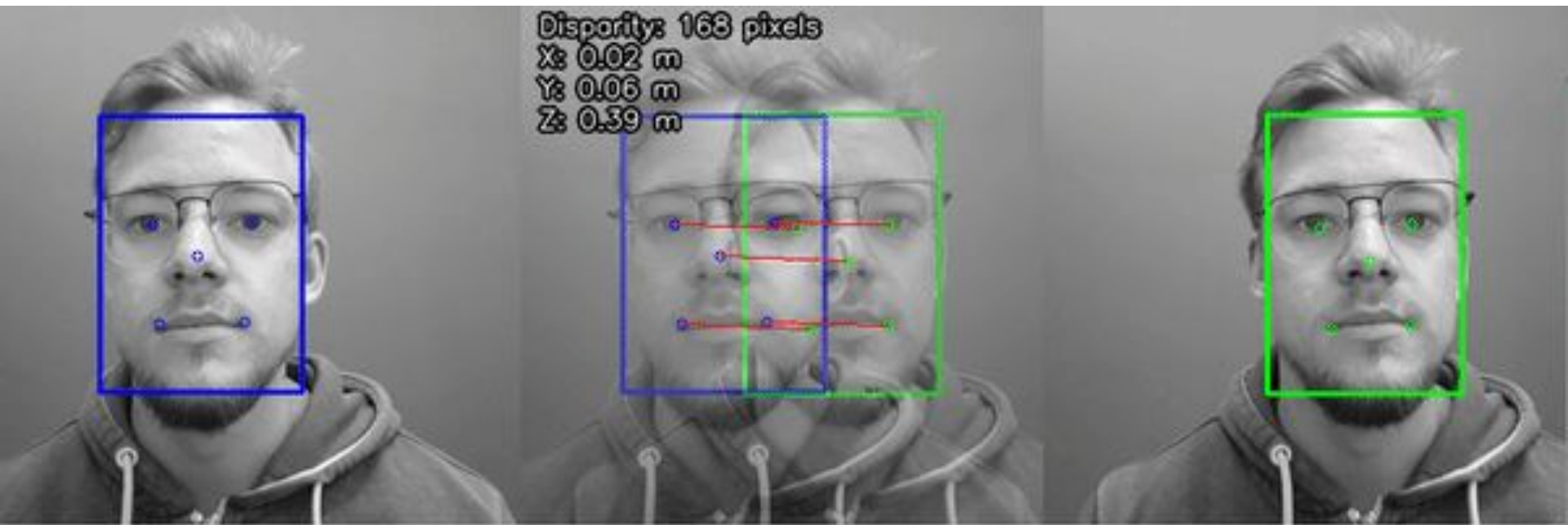




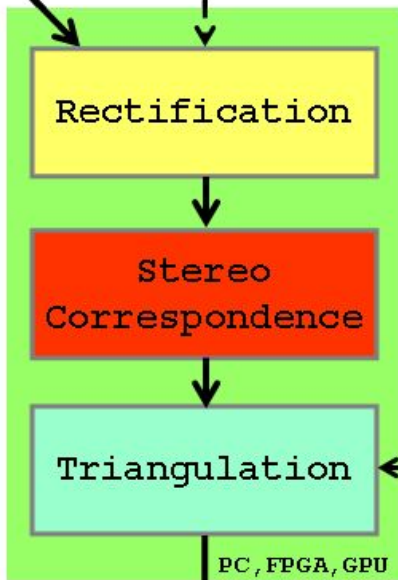
La disparidad **disminuye** con la distancia



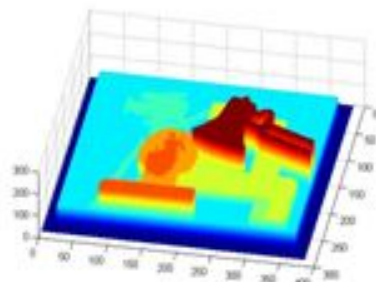
# La disparidad **disminuye** con la distancia

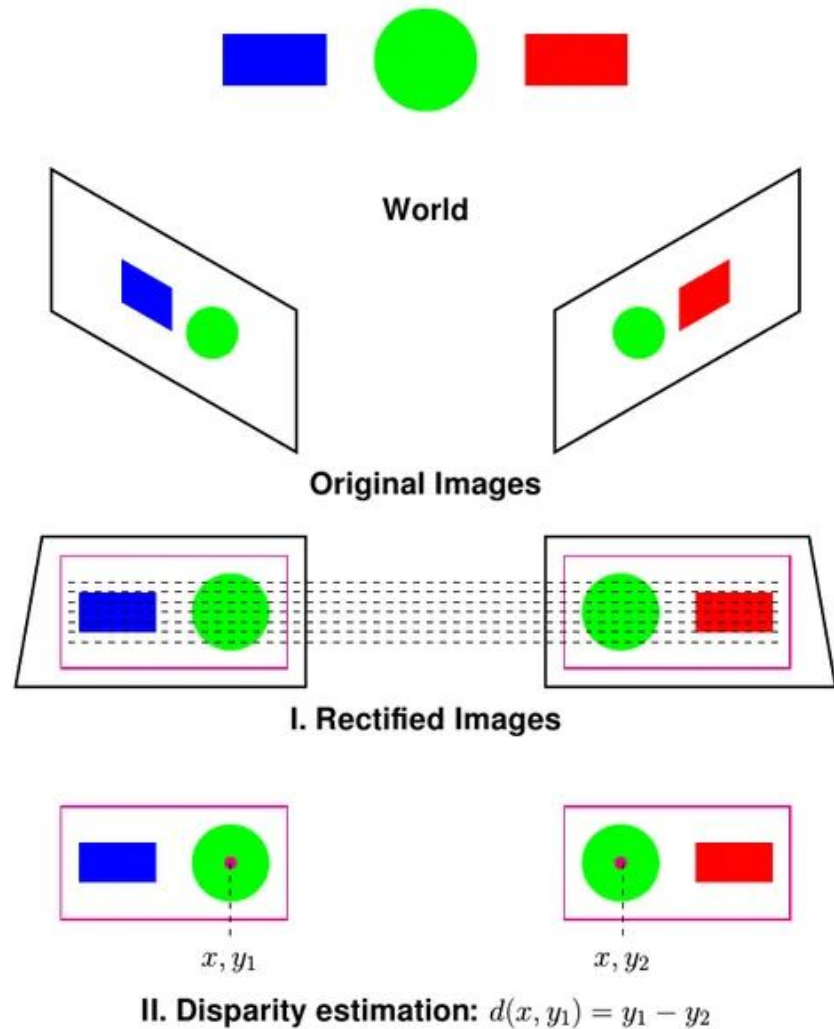
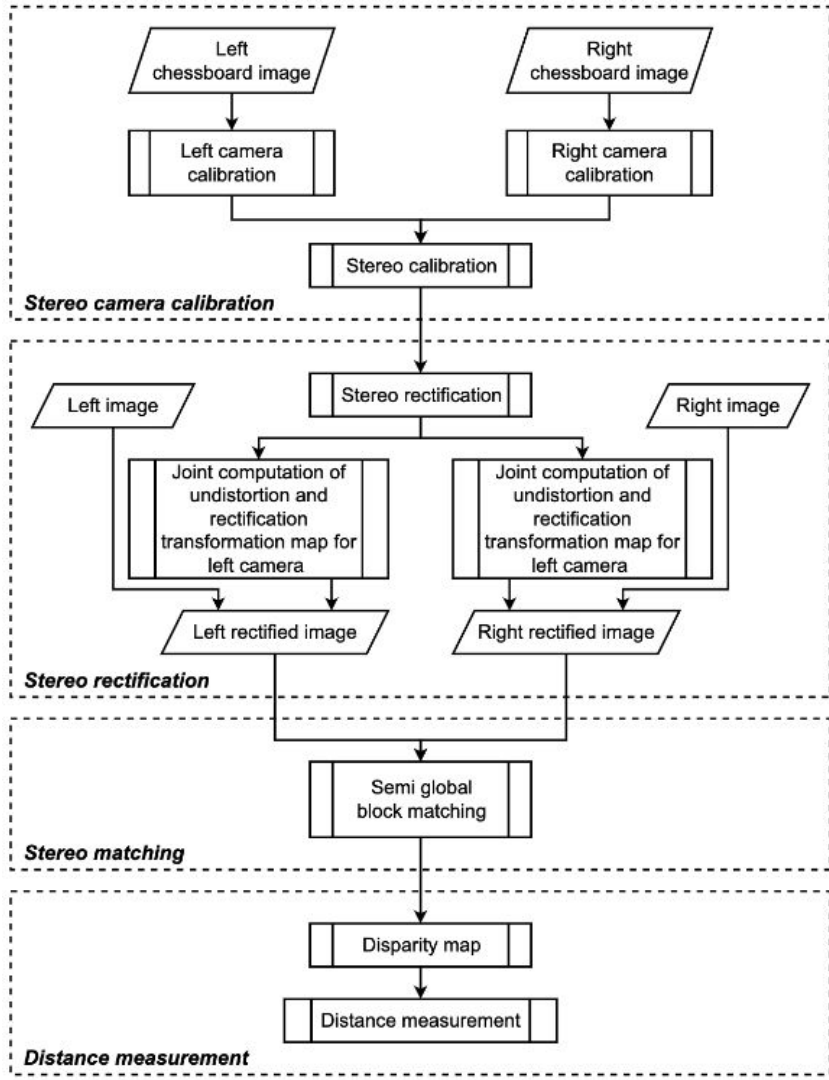


# Passive Stereo (Pipeline)



PC, FPGA, GPU  
3D





¿Cómo obtener la disparidad?

Proyección  
Perspectiva



# Proyección Perspectiva



A group of about ten people, mostly older adults, are gathered around a mini-golf course. Two vertical blue lines are drawn on the image to compare the heights of two women standing in the center. One woman is wearing a white polo shirt and purple pants, while the other is wearing a white polo shirt and dark pants. The background shows a sunny outdoor setting with a fence and some greenery.

**Who is taller?**

A close-up view of a mini-golf green. Three balls are on the green: a black ball on the left, a white ball in the middle, and a red ball on the right. Two red lines originate from the white ball, one pointing towards the black ball and the other pointing towards the red ball, to compare their distances.

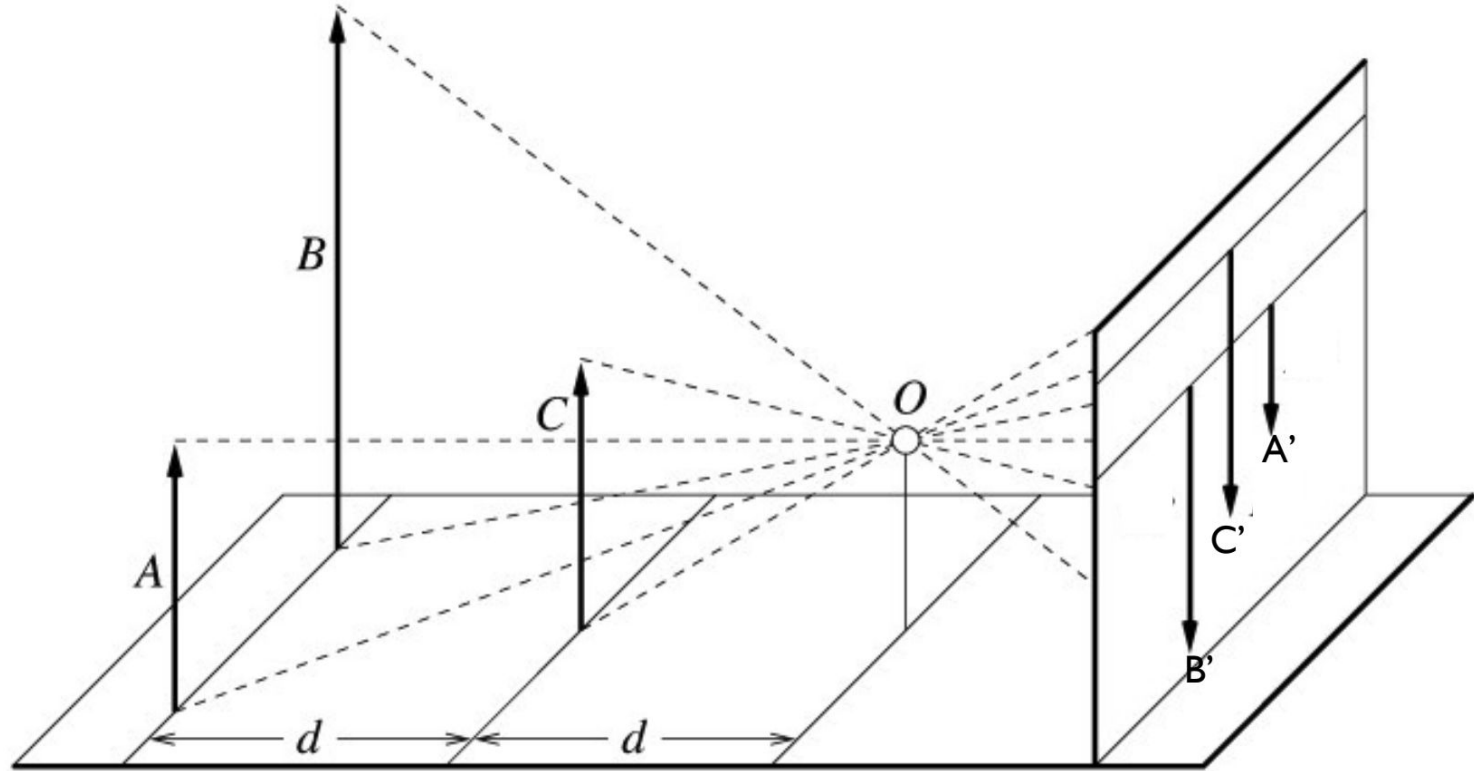
**Which is closer?**





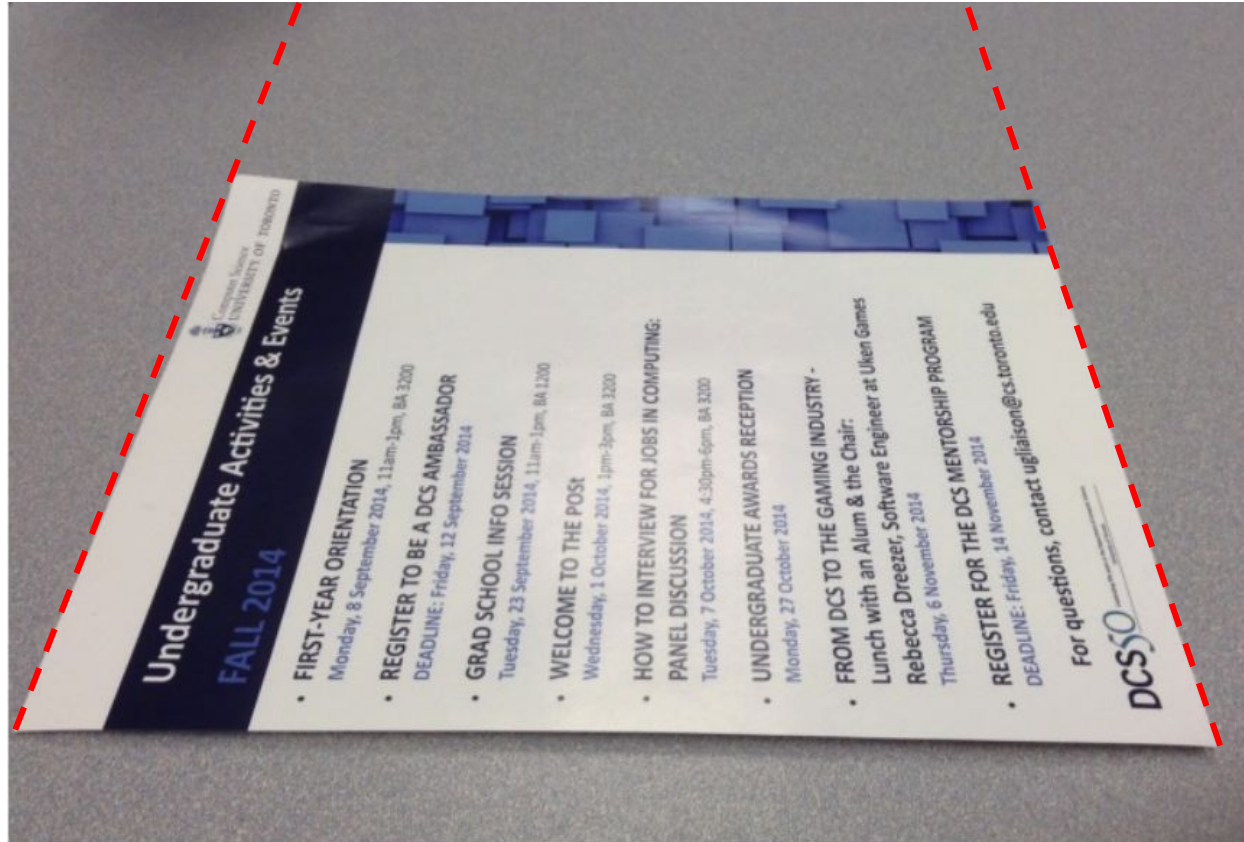


# Geometría de Proyección Perspectiva



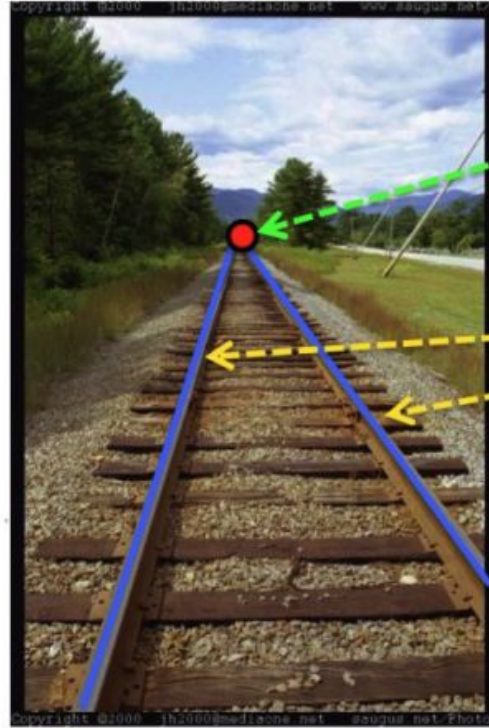
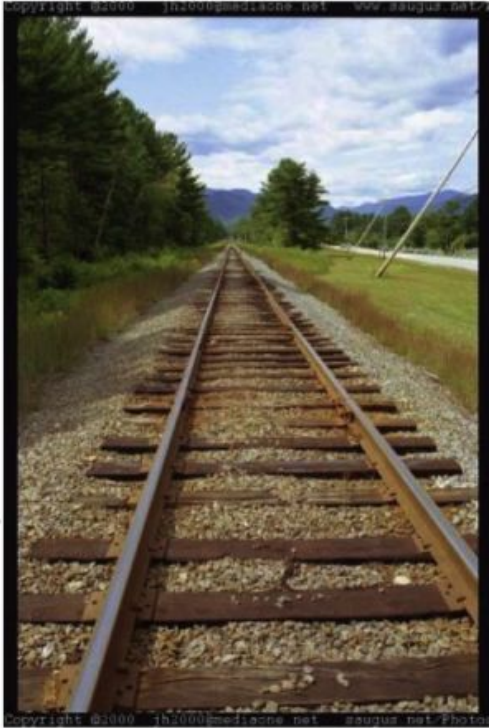
La longitud y el área no se conservan.  $A$  y  $C$  son del mismo tamaño, pero  $A$  se ve más pequeño en el sensor. De manera similar,  $B$  es más grande que  $C$ , pero aparecen del mismo tamaño en el sensor.

# Vanishing Points



Líneas paralelas en el mundo 3D convergen a un punto en 2D (vanishing point)

# Vanishing Points

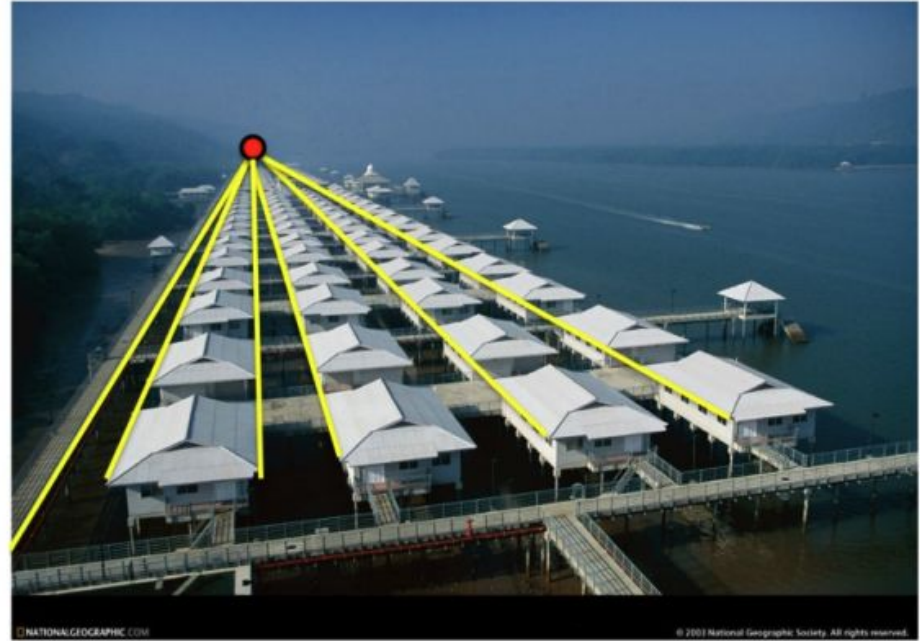


vanishing point

lines parallel in  
the 3D world

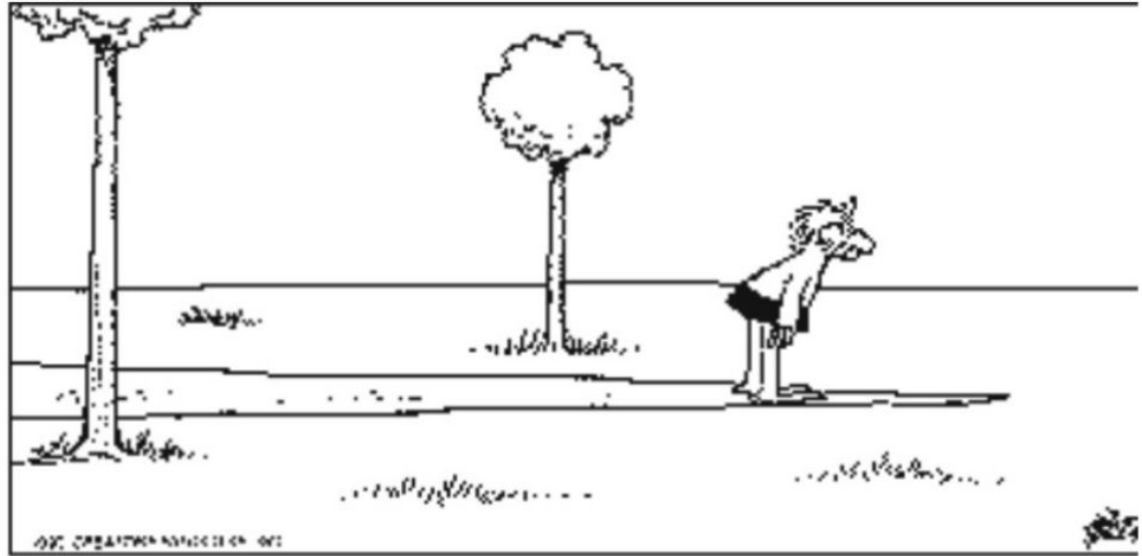
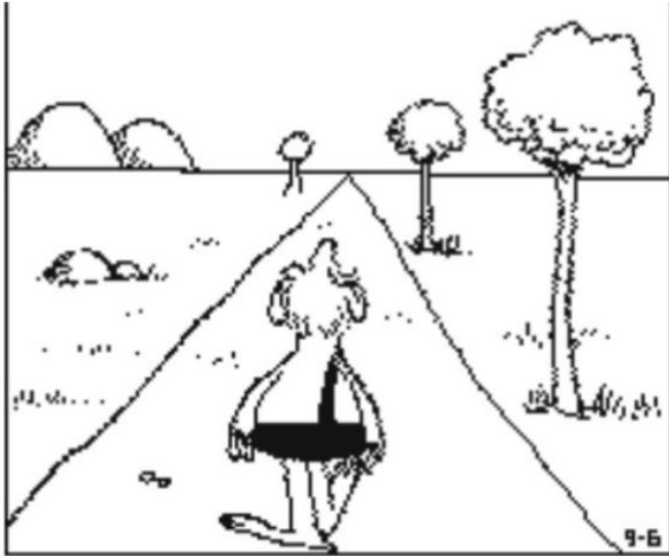
Cada ángulo de dirección diferente en el mundo tiene su propio “vanishing point”

# Vanishing Points



Todas las líneas con la misma dirección en el mundo 3D se intersectan en el mismo vanishing point

# Vanishing Points



**Lo contrario no siempre es verdad:** Líneas que se intersectan en 2D no siempre corresponden a líneas paralelas en 3D

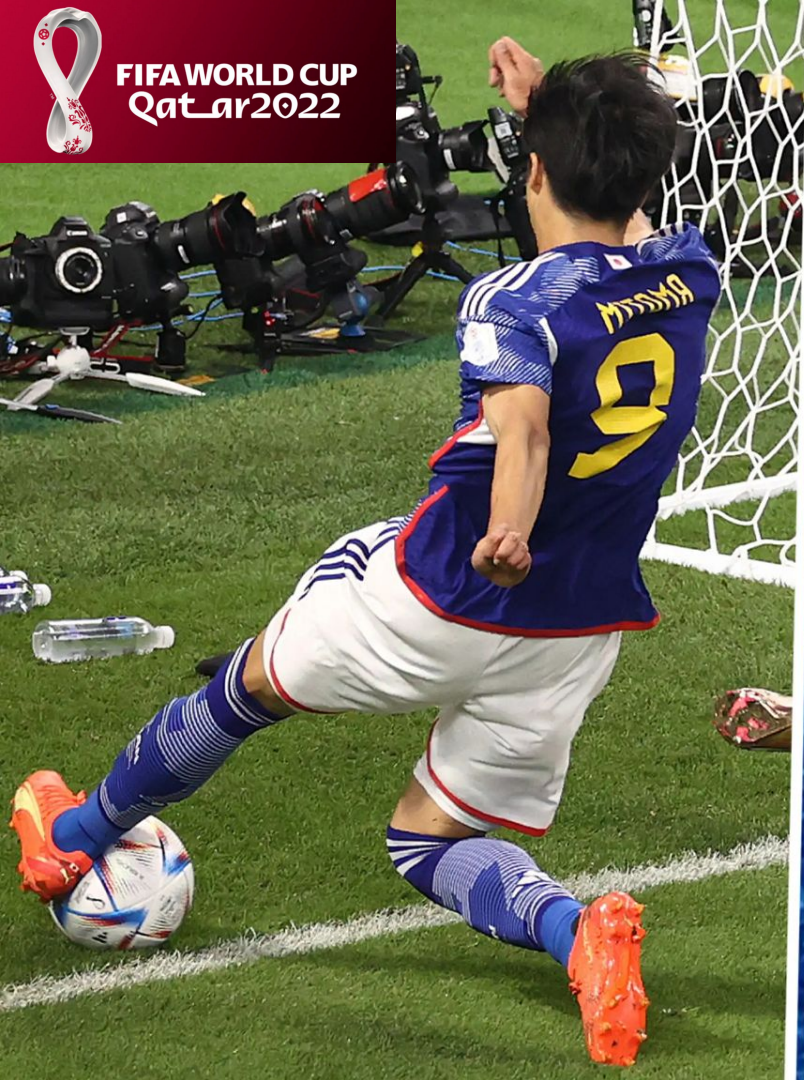
Inter



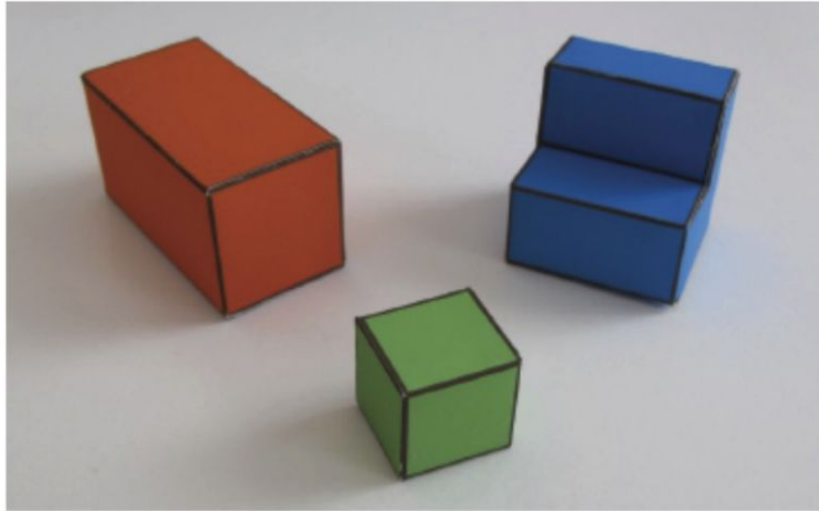




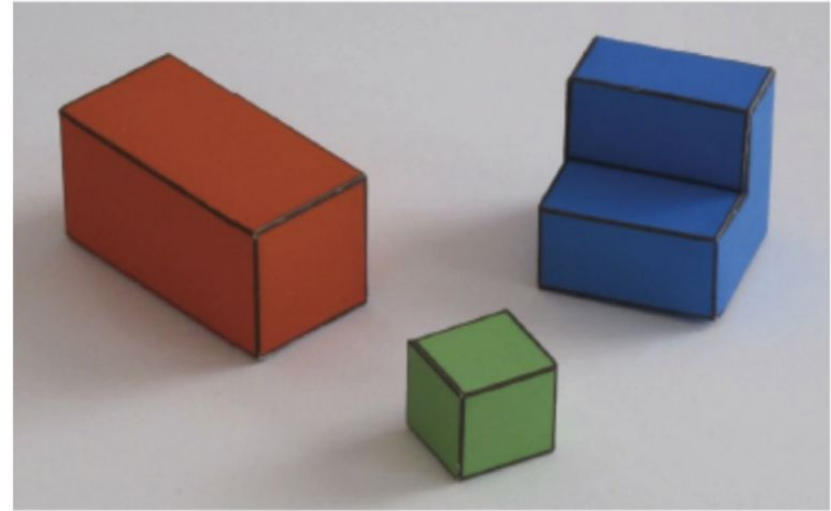
FIFA WORLD CUP  
Qatar 2022



# Tipos de Proyecciones: Perspectiva vs Ortográfica



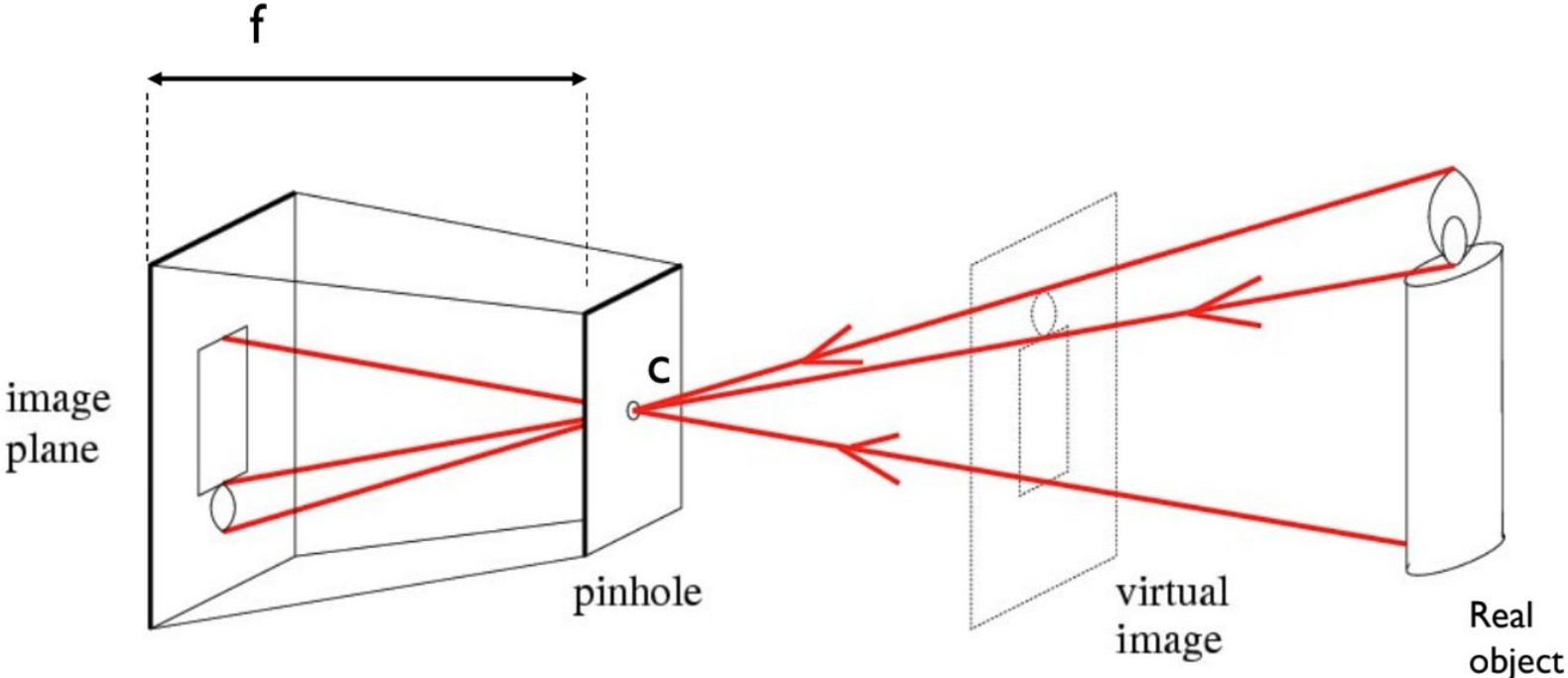
Perspective projection



Parallel (orthographic) projection

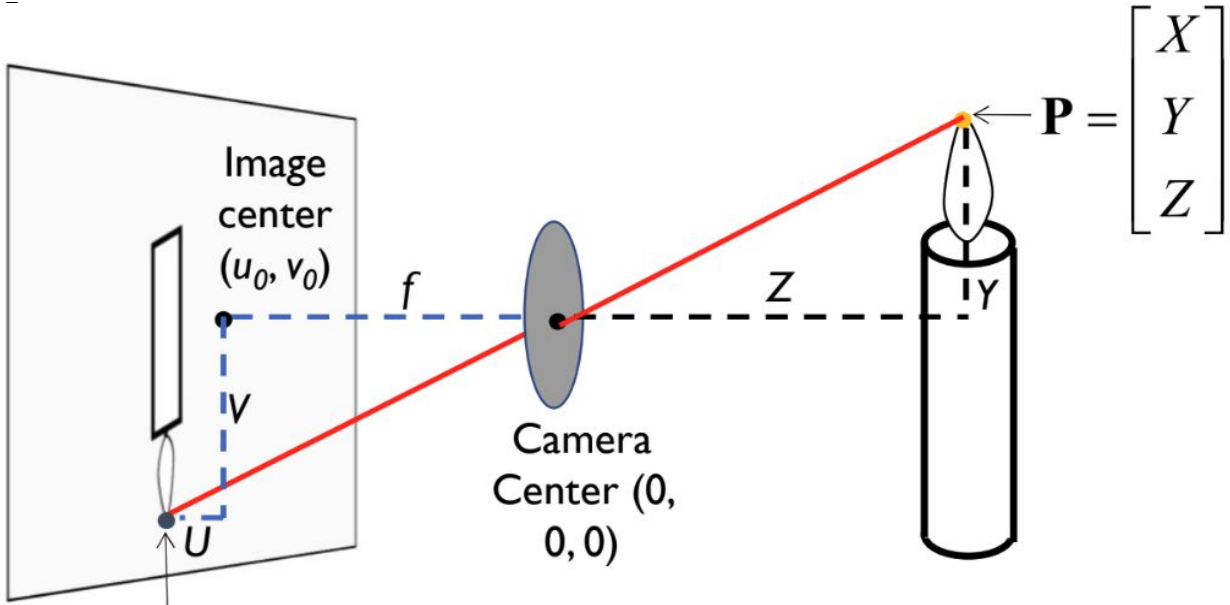
- For perspective projection lines parallel in 3D **are not** parallel in the image.
- For orthographic projection lines parallel in 3D **are** parallel in the image.

# Pinhole Camera Model

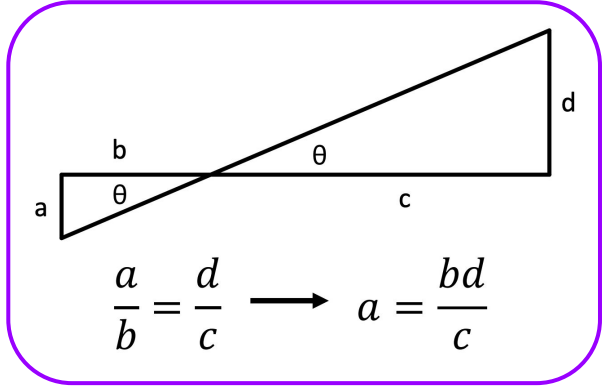


$f$  = Focal length  
 $c$  = Optical center of the camera

# Geometría de Proyección Perspectiva



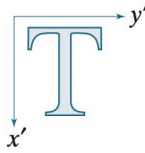
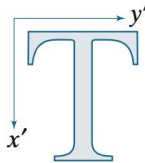
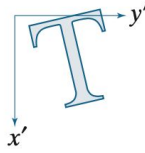
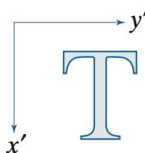
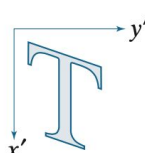
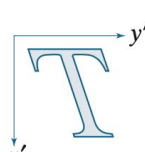
Recordando triángulos



$\mathbf{p} = \begin{bmatrix} U \\ V \end{bmatrix}$   
 $\mathbf{p}$  = distance from image center

$$U = -X * \frac{f}{Z}$$

$$V = -Y * \frac{f}{Z}$$

Transformation Name	Affine Matrix, A	Coordinate Equations	Example
Identity	$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$	$\begin{aligned} x' &= x \\ y' &= y \end{aligned}$	
Scaling/Reflection (For reflection, set one scaling factor to -1 and the other to 0)	$\begin{bmatrix} c_x & 0 & 0 \\ 0 & c_y & 0 \\ 0 & 0 & 1 \end{bmatrix}$	$\begin{aligned} x' &= c_x x \\ y' &= c_y y \end{aligned}$	
Rotation (about the origin)	$\begin{bmatrix} \cos \theta & -\sin \theta & 0 \\ \sin \theta & \cos \theta & 0 \\ 0 & 0 & 1 \end{bmatrix}$	$\begin{aligned} x' &= x \cos \theta - y \sin \theta \\ y' &= x \sin \theta + y \cos \theta \end{aligned}$	
Translation	$\begin{bmatrix} 1 & 0 & t_x \\ 0 & 1 & t_y \\ 0 & 0 & 1 \end{bmatrix}$	$\begin{aligned} x' &= x + t_x \\ y' &= y + t_y \end{aligned}$	
Shear (vertical)	$\begin{bmatrix} 1 & s_v & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$	$\begin{aligned} x' &= x + s_v y \\ y' &= y \end{aligned}$	
Shear (horizontal)	$\begin{bmatrix} 1 & 0 & 0 \\ s_h & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$	$\begin{aligned} x' &= x \\ y' &= s_h x + y \end{aligned}$	

# Geometría de Proyección Perspectiva: Ecuaciones

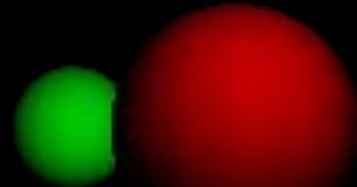
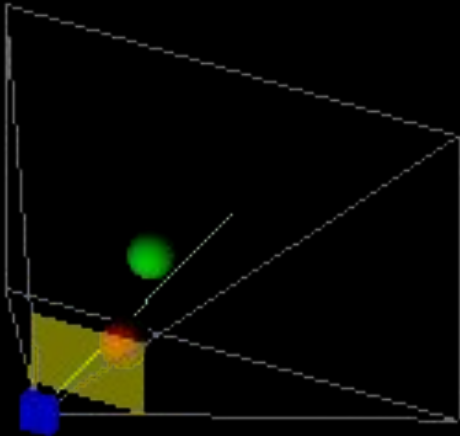
Una cámara y su posición en el espacio está descrita por varios parámetros:

- **Translación (T)** del centro óptico de la imagen al de las coordenadas reales
- **Rotación (R)** del plano imagen
- **Distancia focal (f), punto principal ( $u_0, v_0$ ), tamaño de pixel ( $s_x, s_y$ )**
- **Azul = Extrínsecos; Rojo = Intrínsecos**

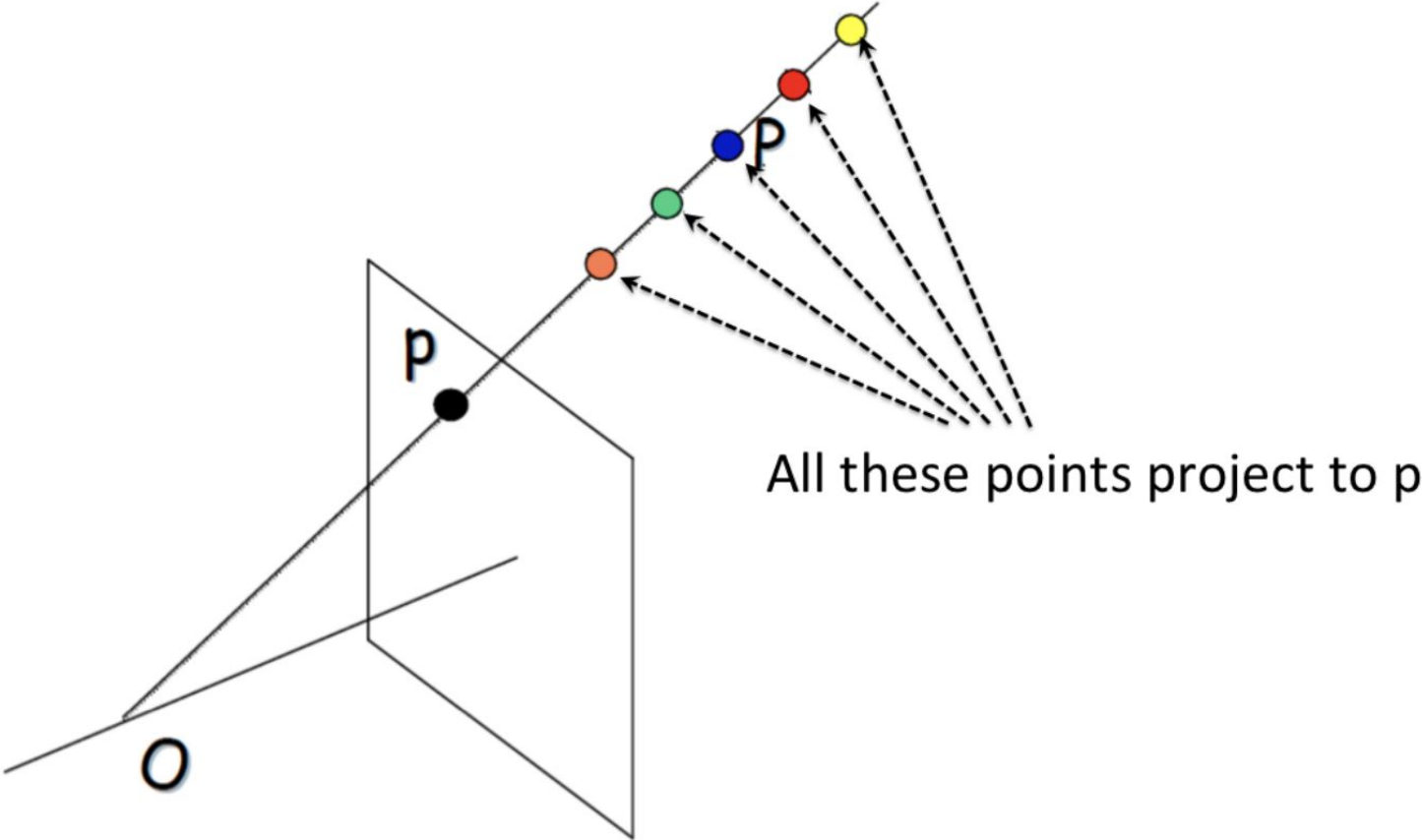
$$\mathbf{p} = \mathbf{K}[\mathbf{R} \quad \mathbf{t}]\mathbf{P} \quad \longrightarrow \quad w \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \underbrace{\begin{bmatrix} f_x & s & u_0 \\ 0 & f_y & v_0 \\ 0 & 0 & 1 \end{bmatrix}}_{\text{Intrínsecos}} \underbrace{\begin{bmatrix} r_{11} & r_{12} & r_{13} & t_x \\ r_{21} & r_{22} & r_{23} & t_y \\ r_{31} & r_{32} & r_{33} & t_z \end{bmatrix}}_{\text{Extrínsecos}} \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix}$$

# Jugando con los parámetros Intrínsecos y Extrínsecos

[https://ksimek.github.io/perspective\\_camera\\_toy.html](https://ksimek.github.io/perspective_camera_toy.html)

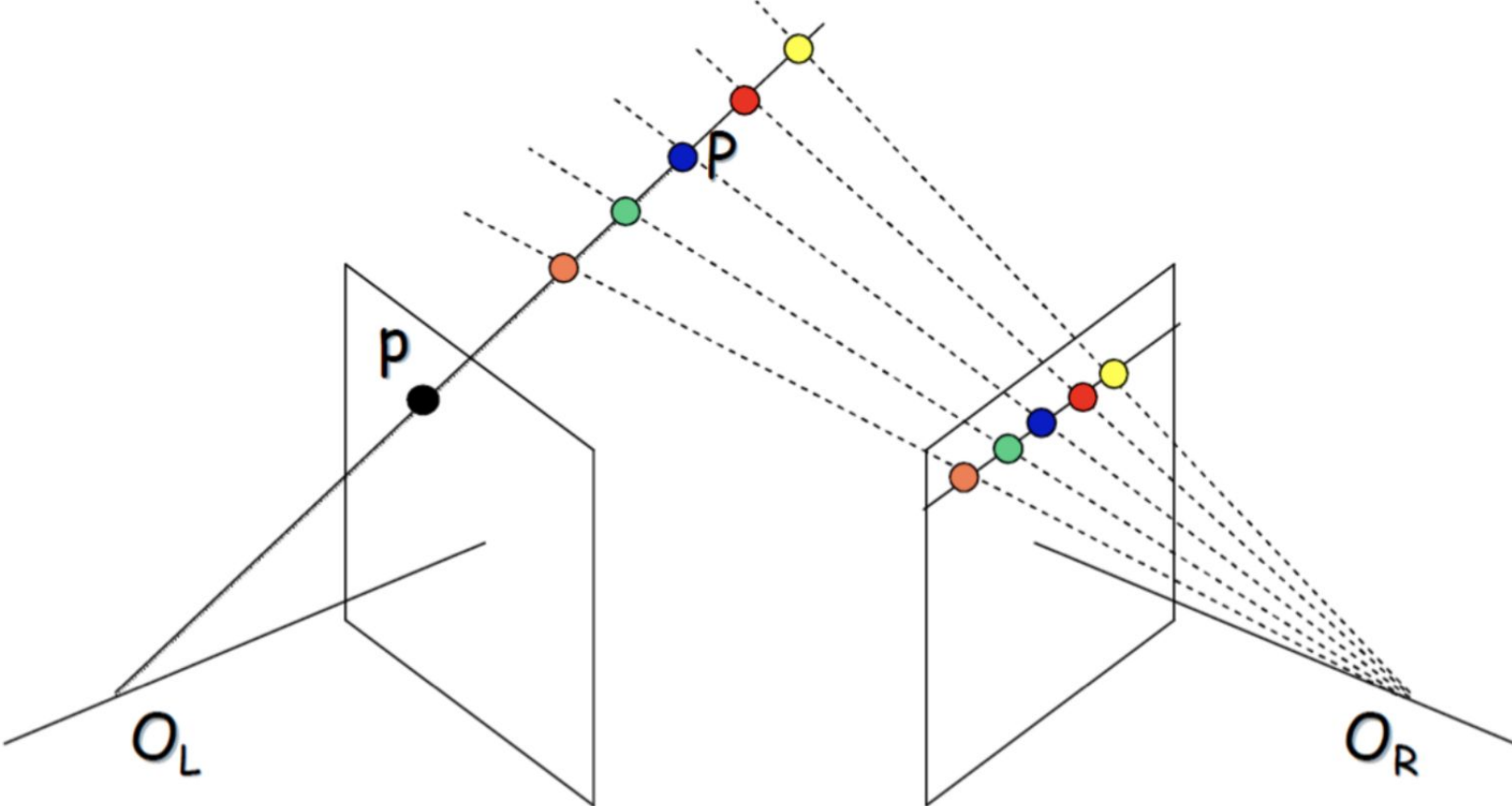


# Geometría de Visión Stereo

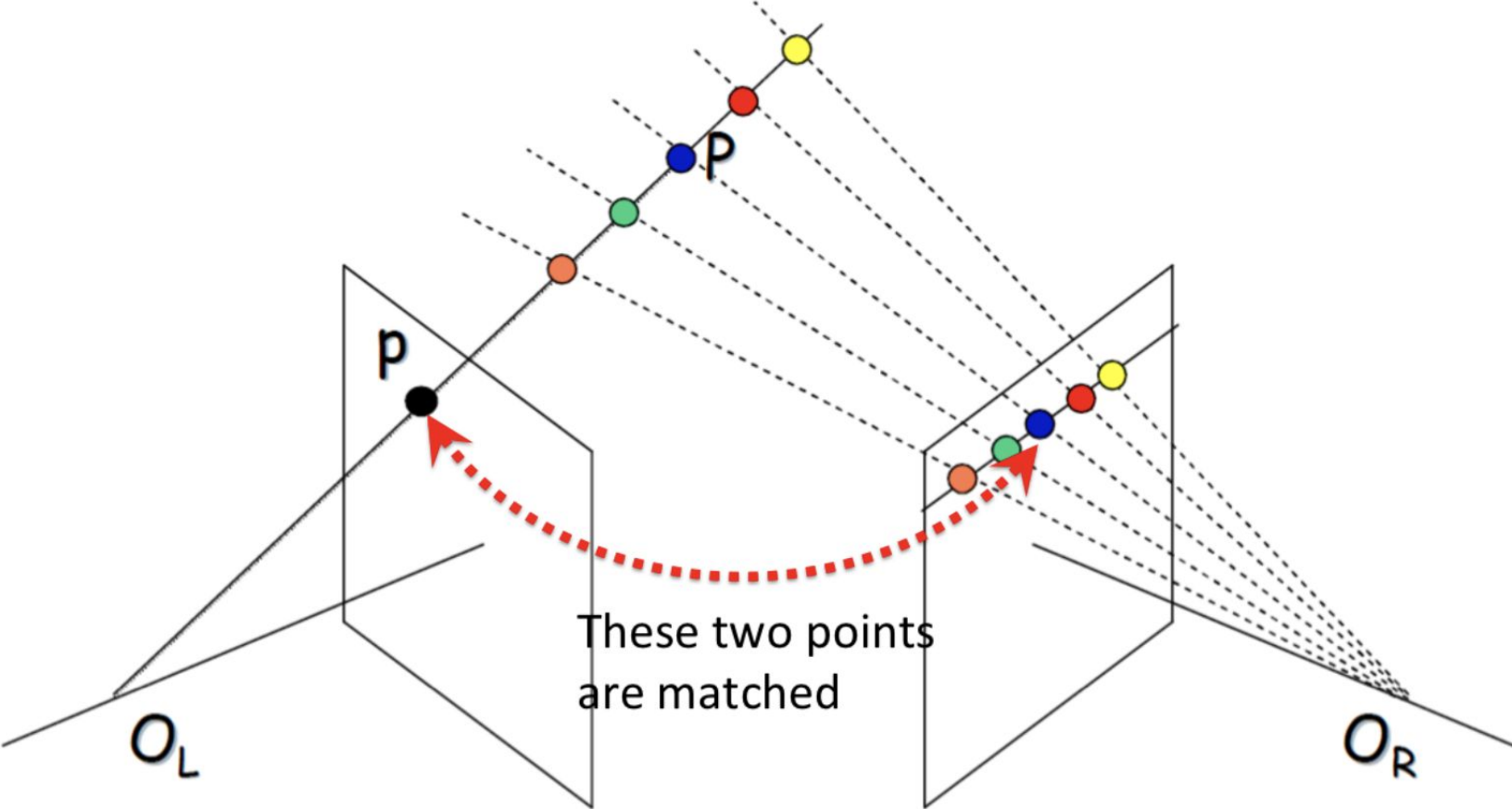




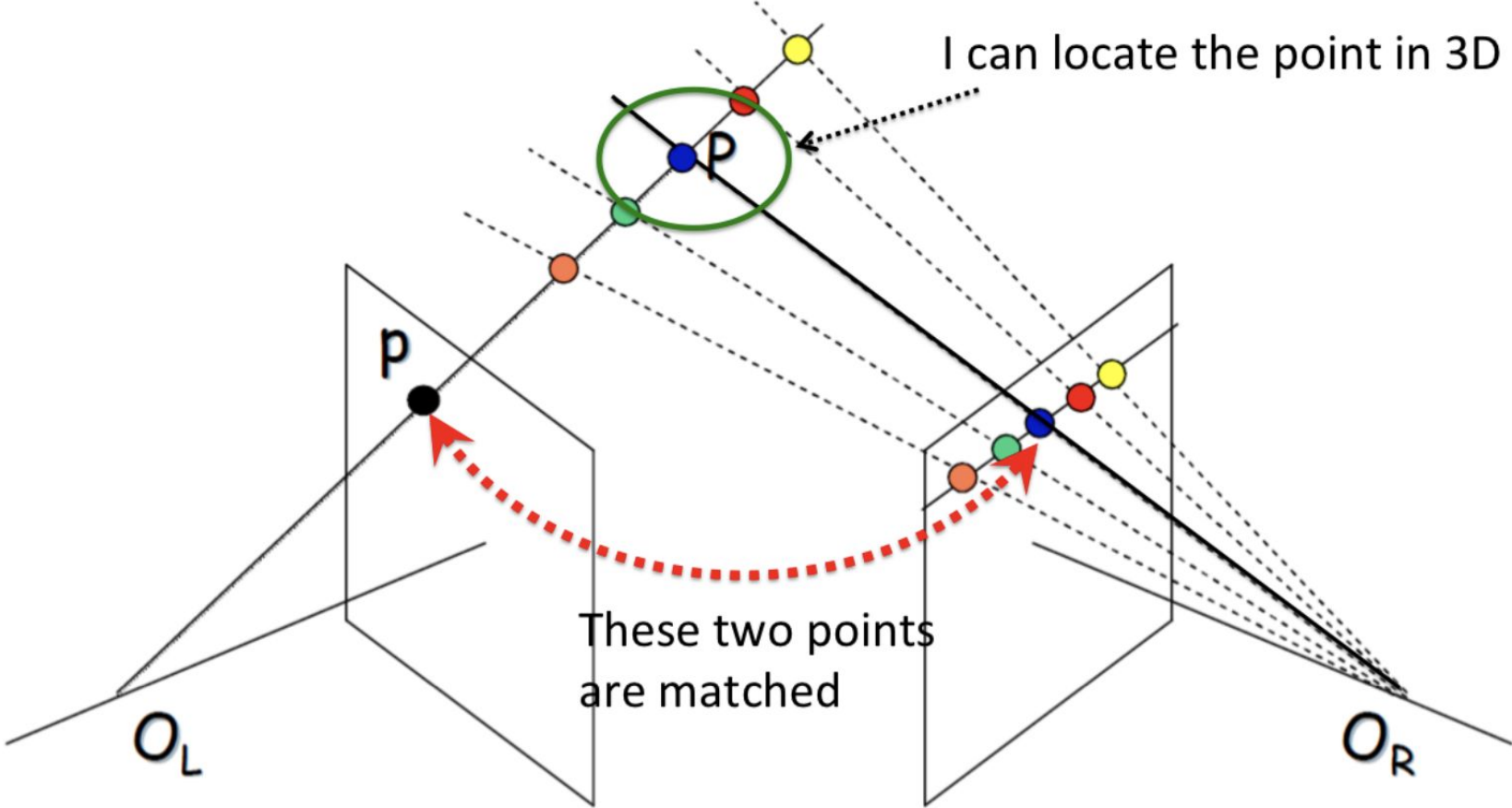
# Geometría de Visión Stereo



# Geometría de Visión Stereo

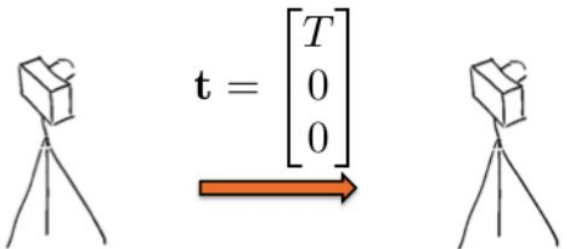
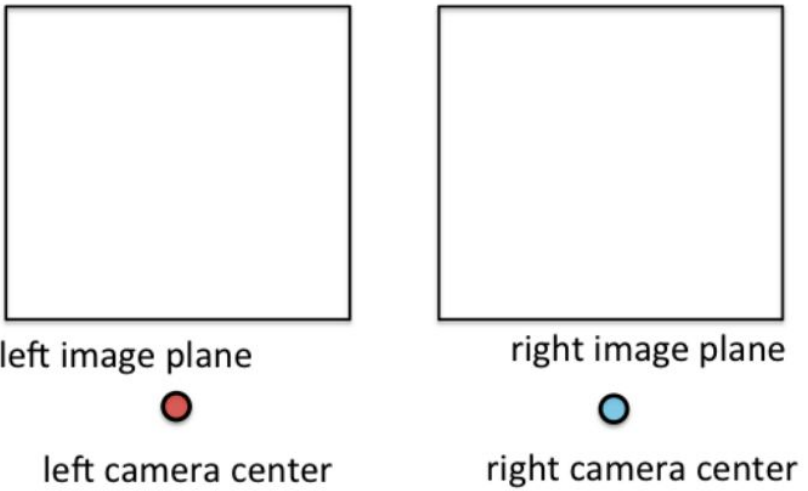


# Geometría de Visión Stereo

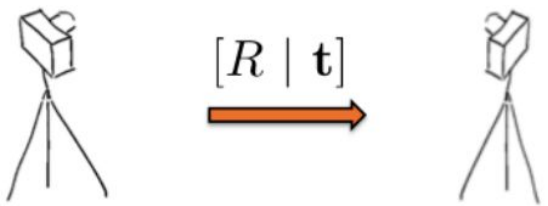
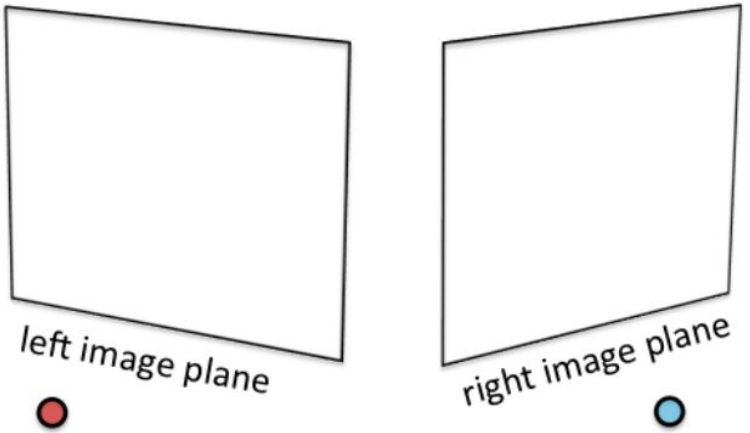


# Geometría de Visión Stereo

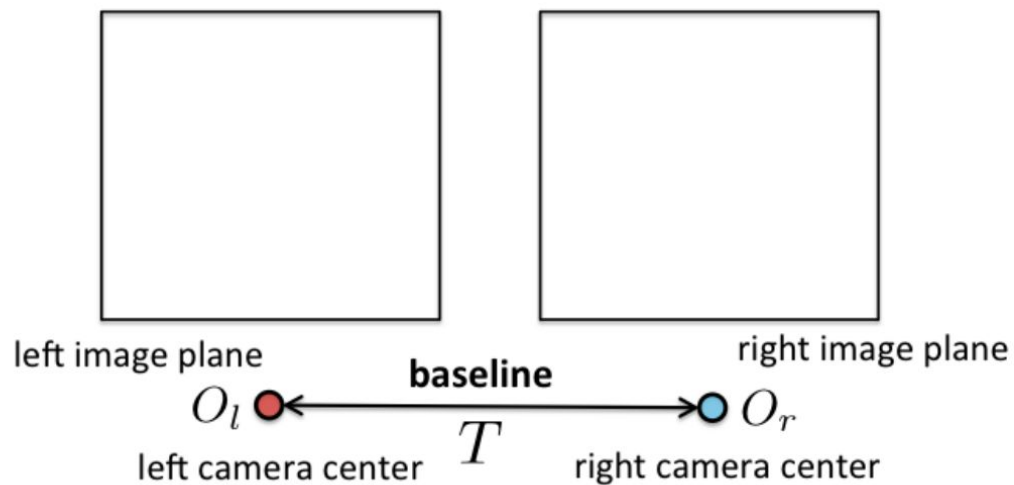
## Paralelas (Caso simple)



## General (Rectificar)



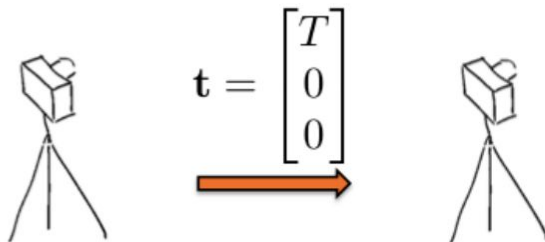
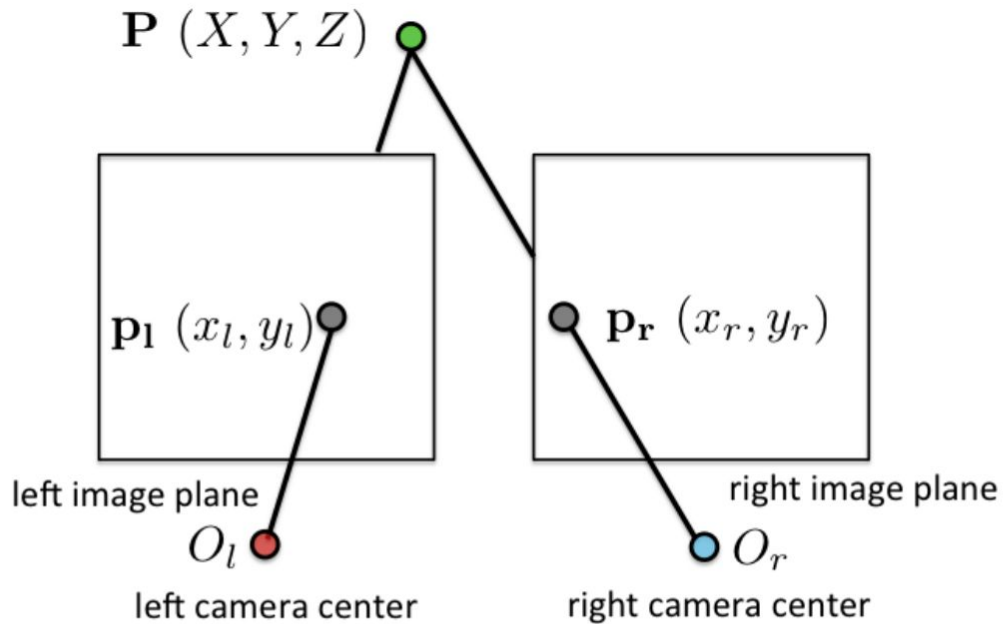
# Caso simple



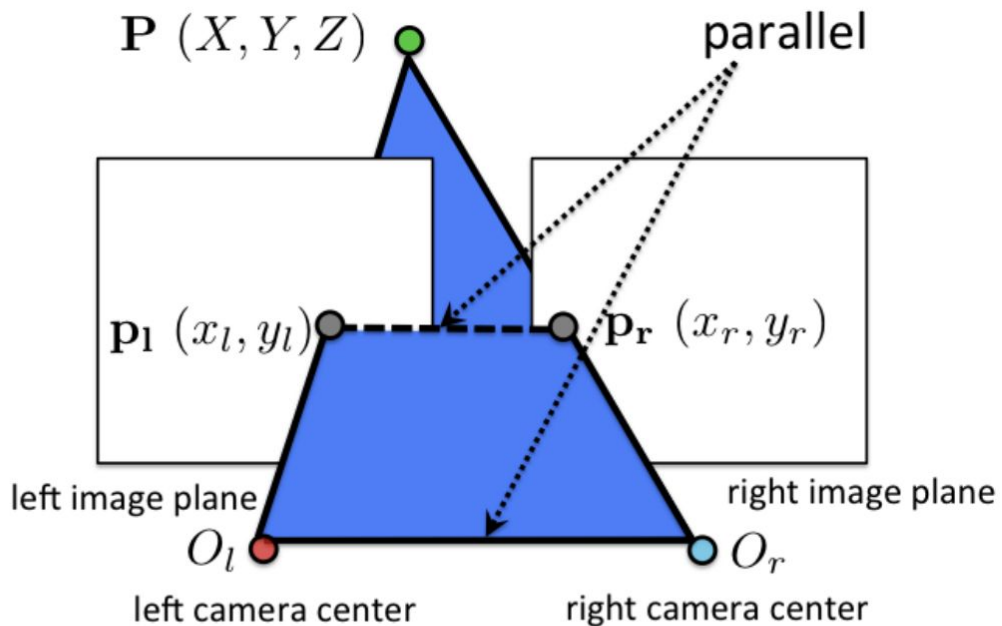
$$\mathbf{t} = \begin{bmatrix} T \\ 0 \\ 0 \end{bmatrix}$$

The right camera  
is shifted to the  
right in X direction

# Caso simple



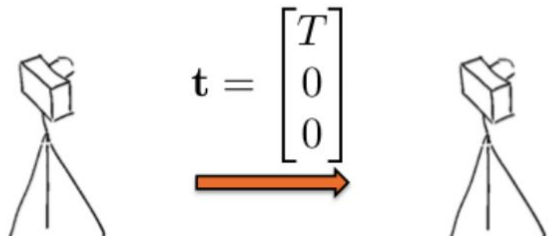
# Caso simple



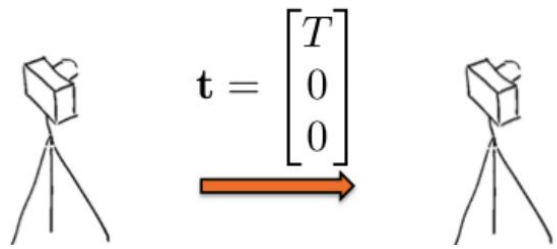
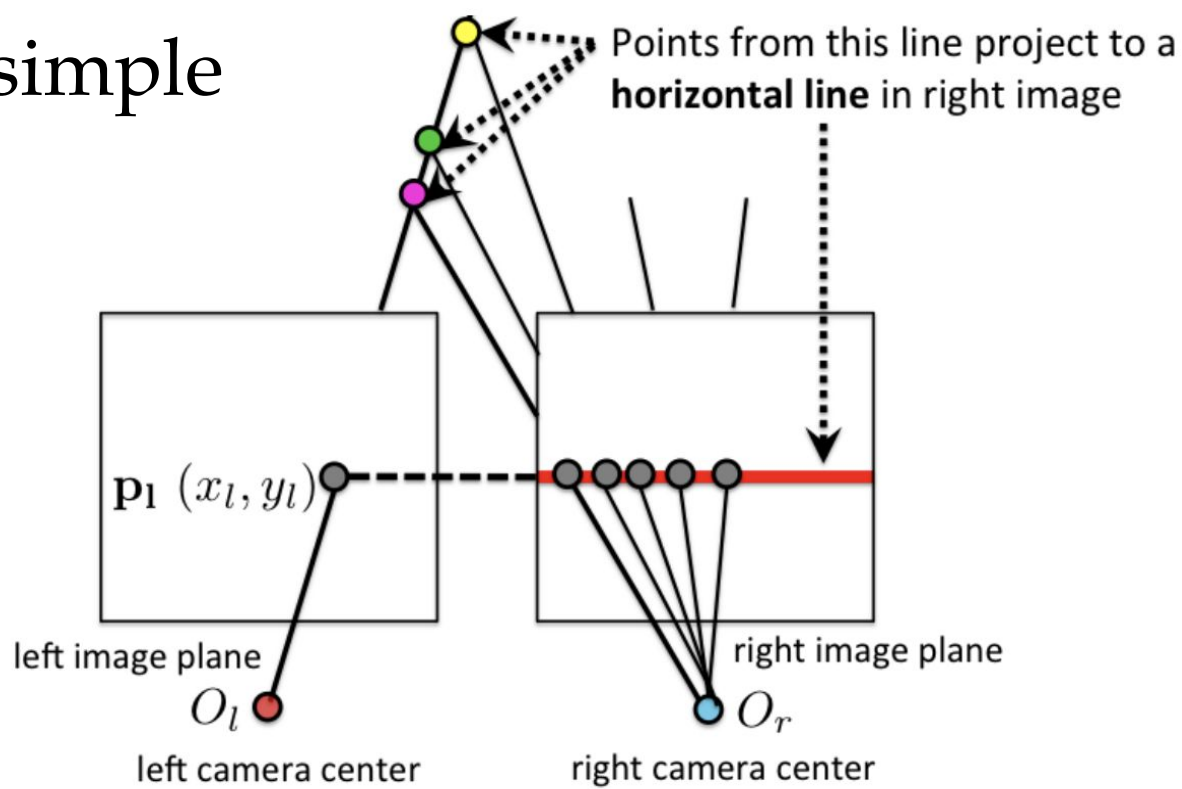
So:  $y_r = y_l$

Nota:

- Olvidémonos de él
- Para poder olvidarnos de él, las cámaras deben estar **rectificadas**

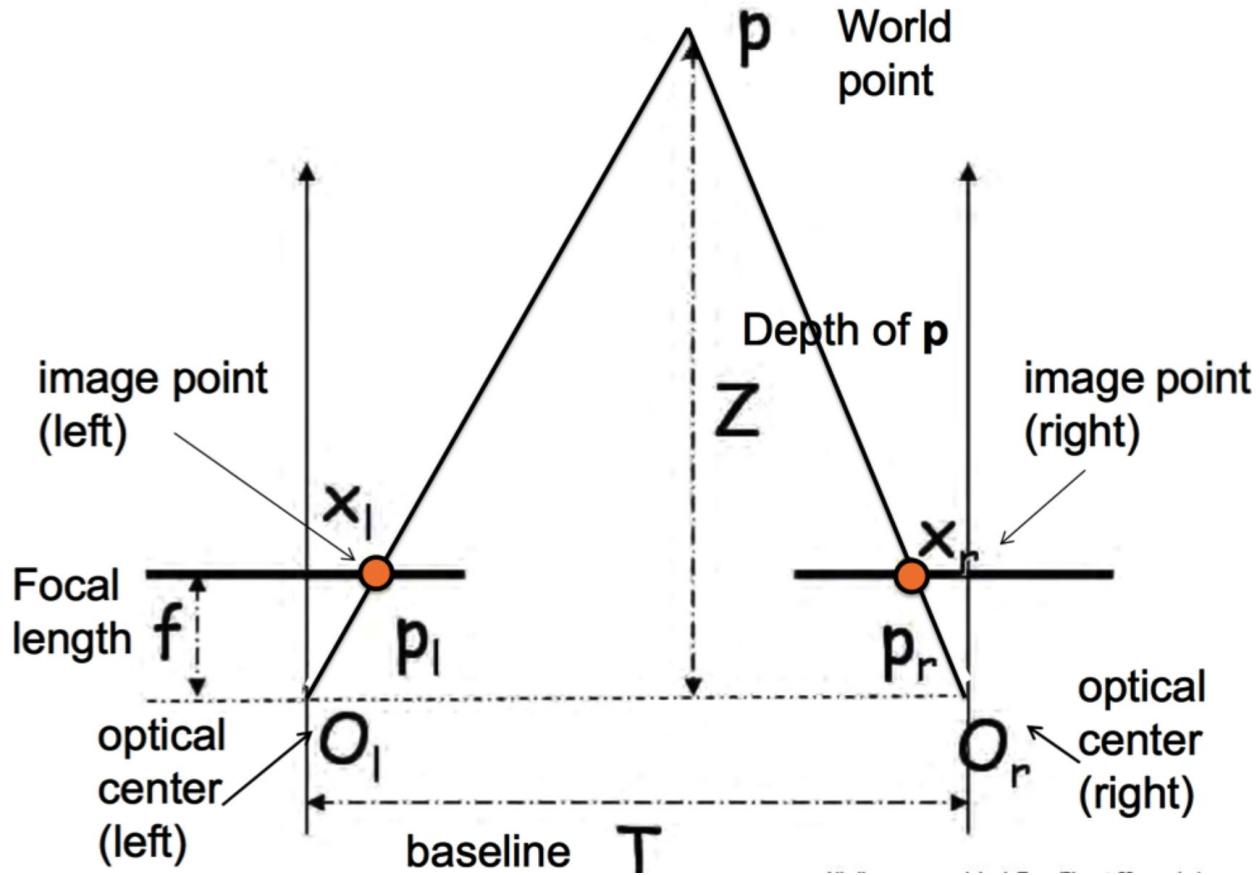


# Caso simple

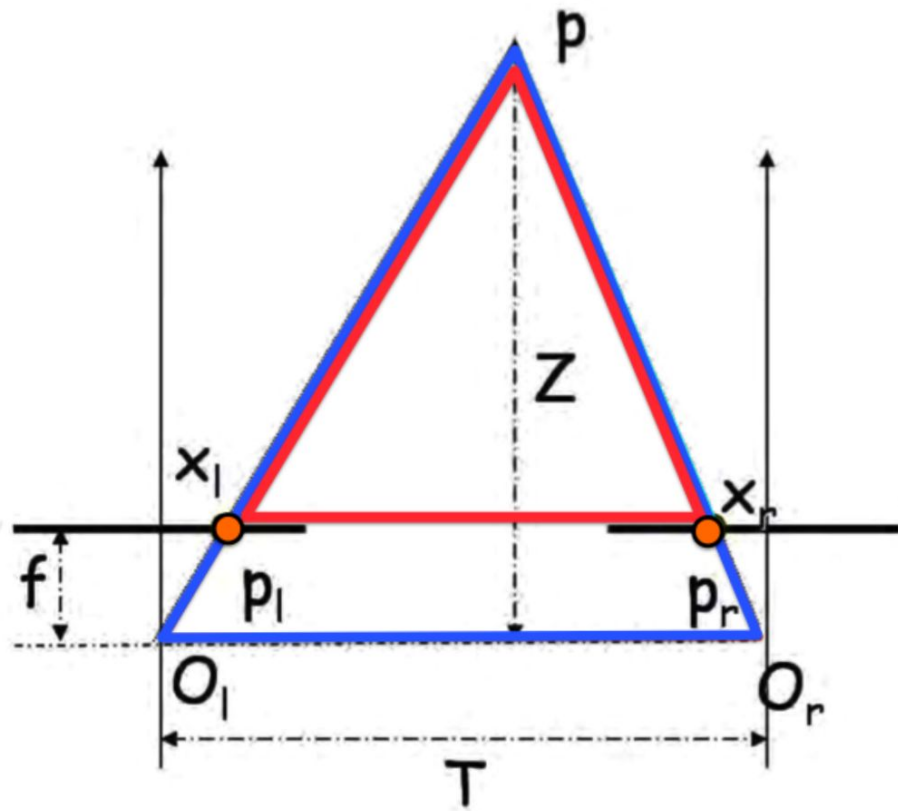




# Calculando Profundidad



# Calculando Profundidad



Similar triangles:

$$\frac{T}{Z} = \frac{T + x_r - x_l}{Z - f}$$

The diagram shows the derivation of the depth equation. The equation  $Z = \frac{f \cdot T}{x_l - x_r}$  is enclosed in an orange rounded rectangle. An orange arrow points from the similar triangles equation above to this equation. Labels with arrows point to the terms in the equation: "focal length" points to  $f$ , "disparity" points to  $x_l - x_r$ , and "baseline" points to  $T$ .

$$Z = \frac{f \cdot T}{x_l - x_r}$$

# Buscando $(x_l, x_r)$



# Buscando $(x_l, x_r)$

$(x_r)$

We are looking for this point



left image

$x_l$



right image

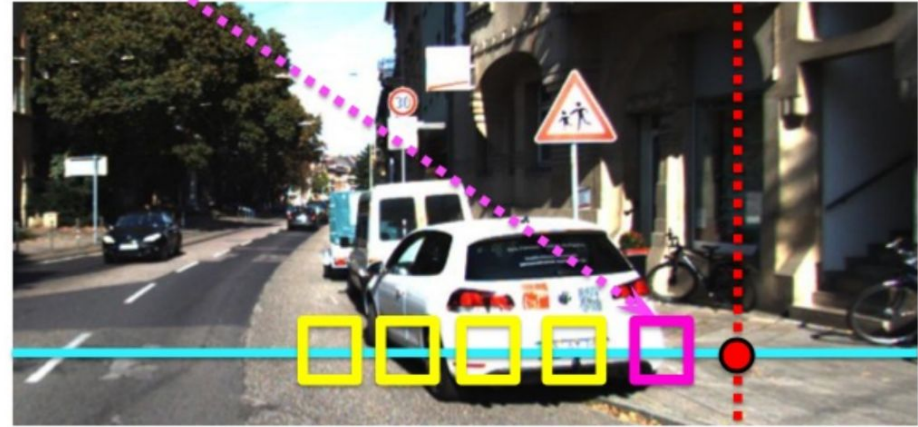
$x_l$

# Buscando $(x_l, x_r)$

Most similar. A match!



left image



right image

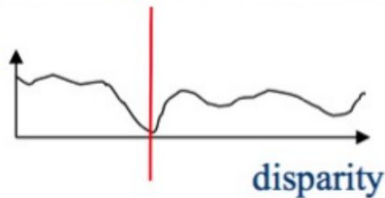
# Buscando $(x_l, x_r)$



left image

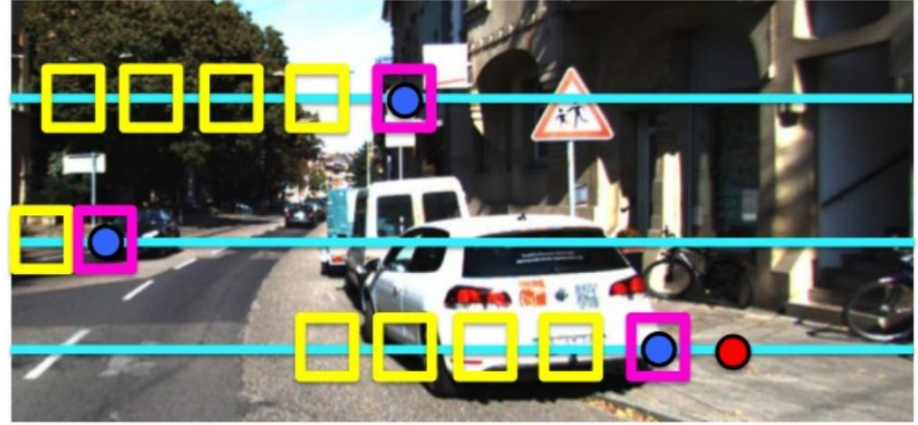


Matching cost

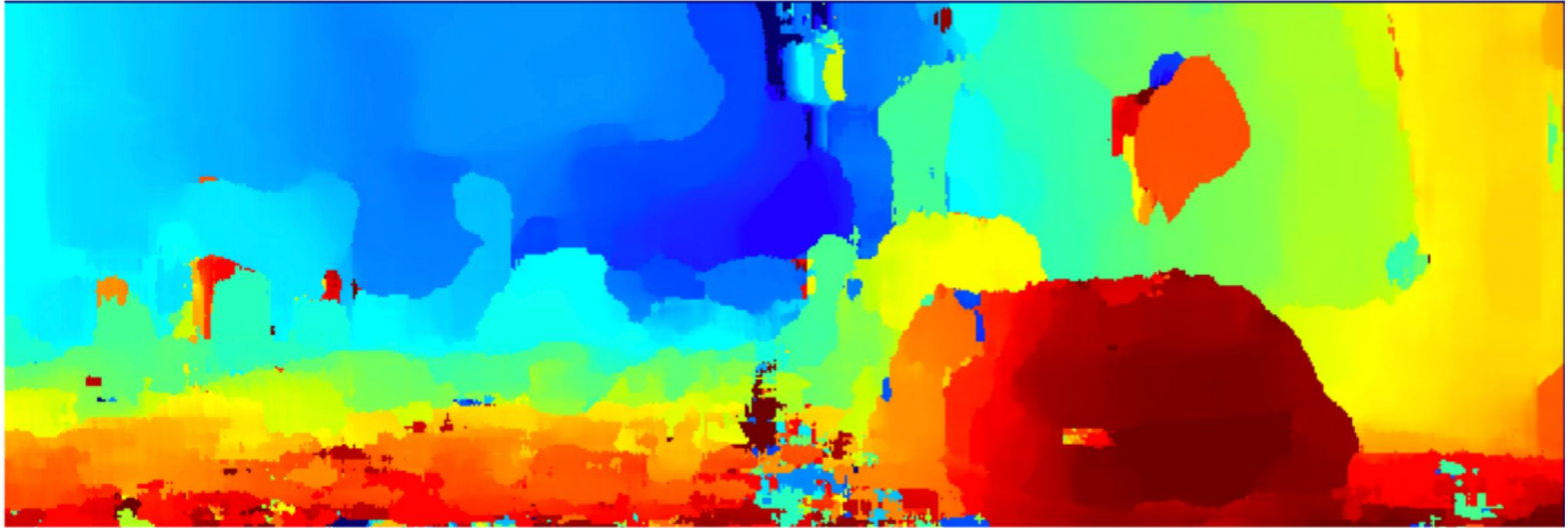


$$SSD(\text{patch}_l, \text{patch}_r) = \sum_x \sum_y (I_{\text{patch}_l}(x, y) - I_{\text{patch}_r}(x, y))^2$$

Buscando  $(x_l, x_r)$ : Para todos los “parches”

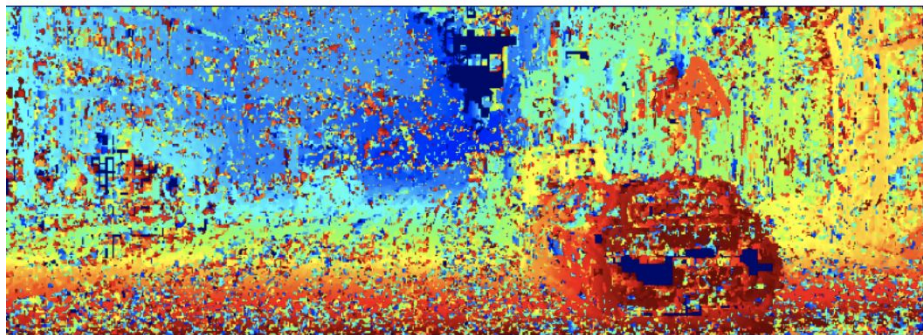


Buscando  $(x_l, x_r)$ : Para todos los “parches”

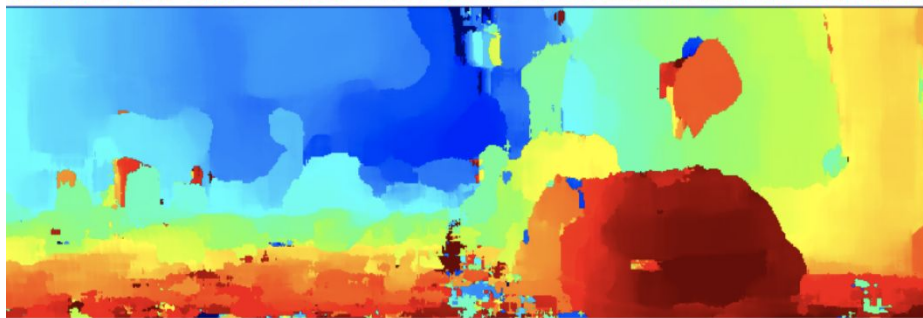




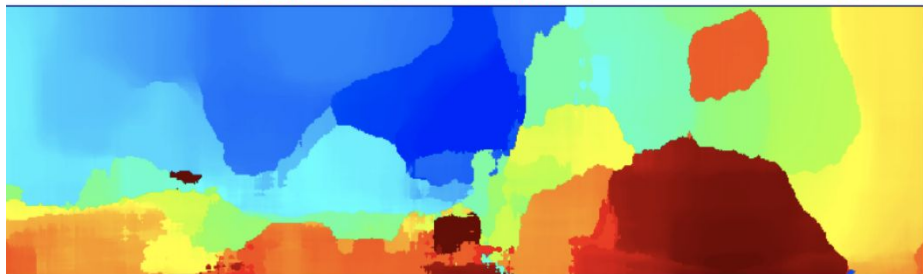
# Disparidad para diferentes tamaños de “parche”



patch size = 5



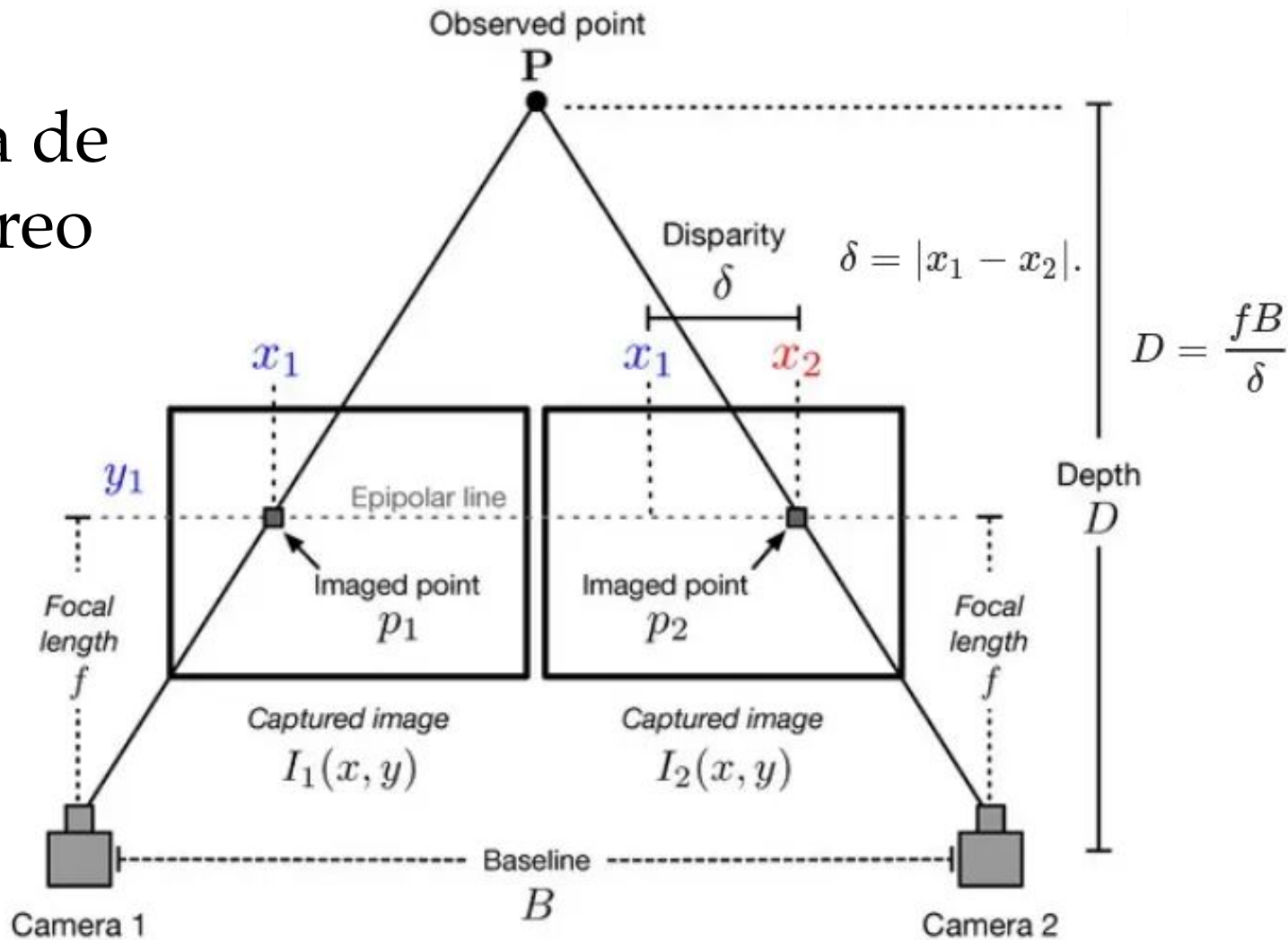
patch size = 35



patch size = 85

**Más grande = menos detalles (ruido)**

# Resumen: Geometría de Visión Stereo



# Demo: Stereo Vision + Tracking

frame left

Distance: 43.3

93%

FPS: 98

frame right

Distance: 43.3

94%

FPS: 98

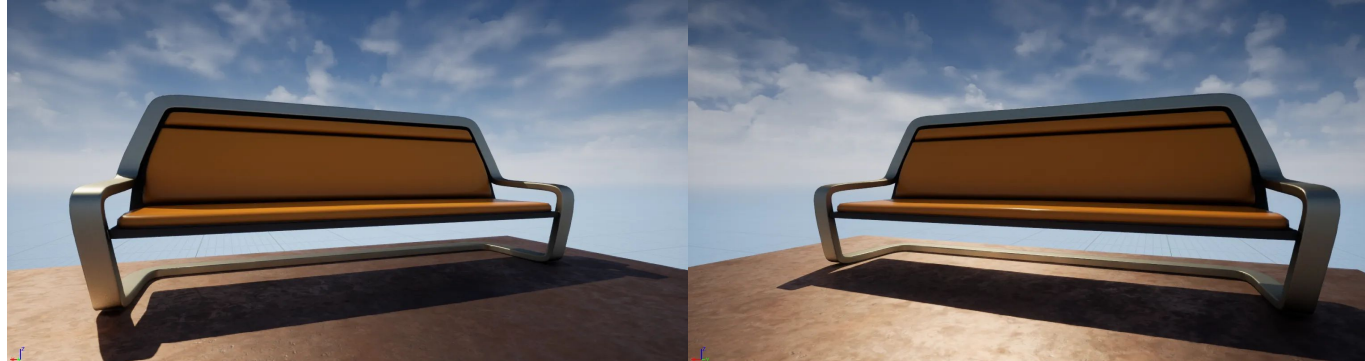
```
42 succes_right, frame_right = cap_right.read()
43 succes_left, frame_left = cap_left.read()
44
```

Run: stereoVision

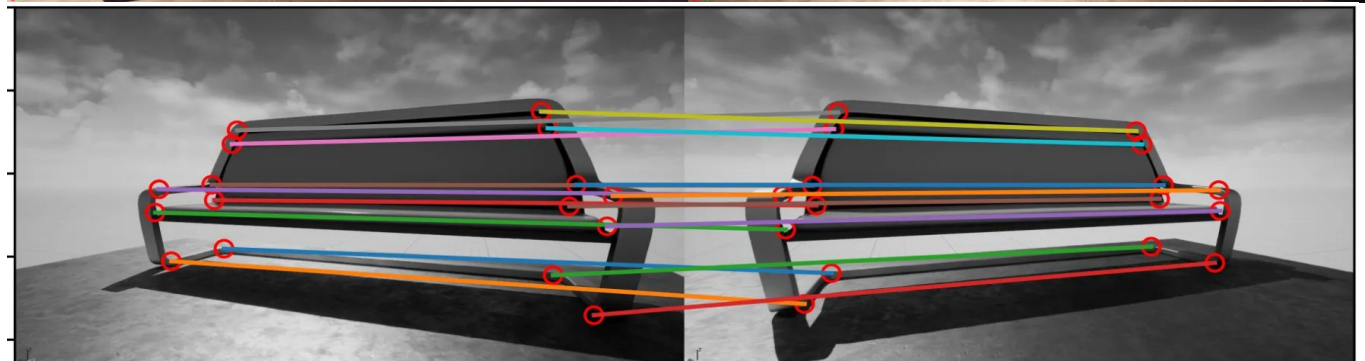
- depth: 44.1
- Depth: 42.5
- Depth: 42.2
- Depth: 41.6
- Depth: 41.8

Caso NO ideal  
(Haciendo rectificaci3n)

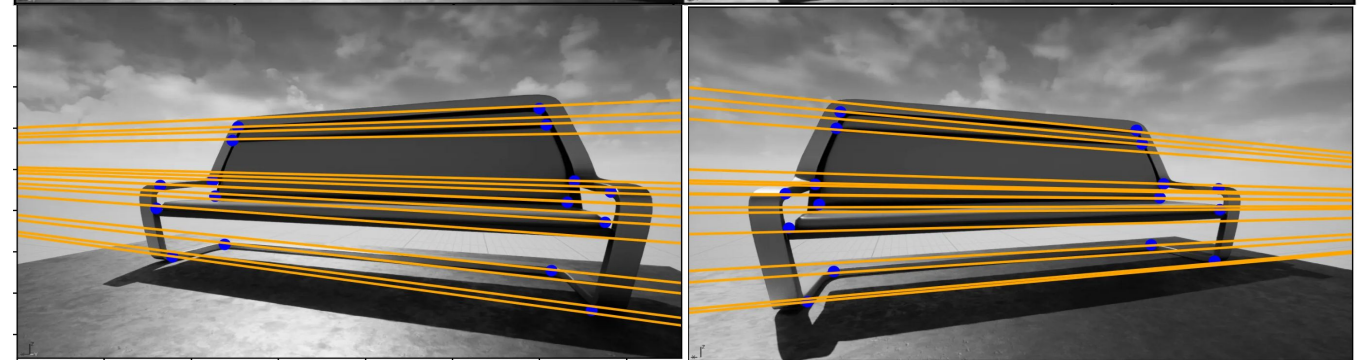
Adquirir  
imágenes



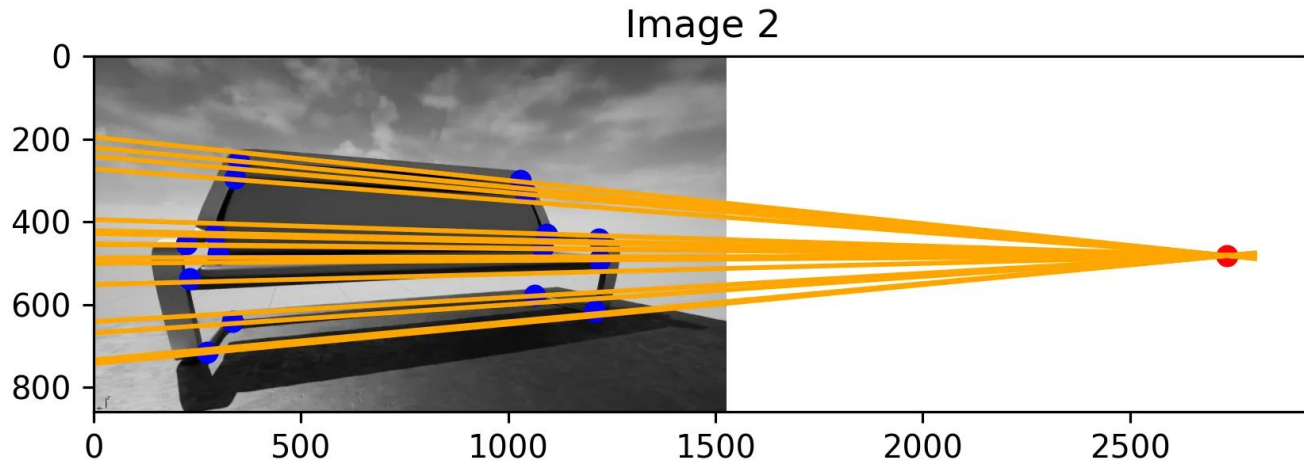
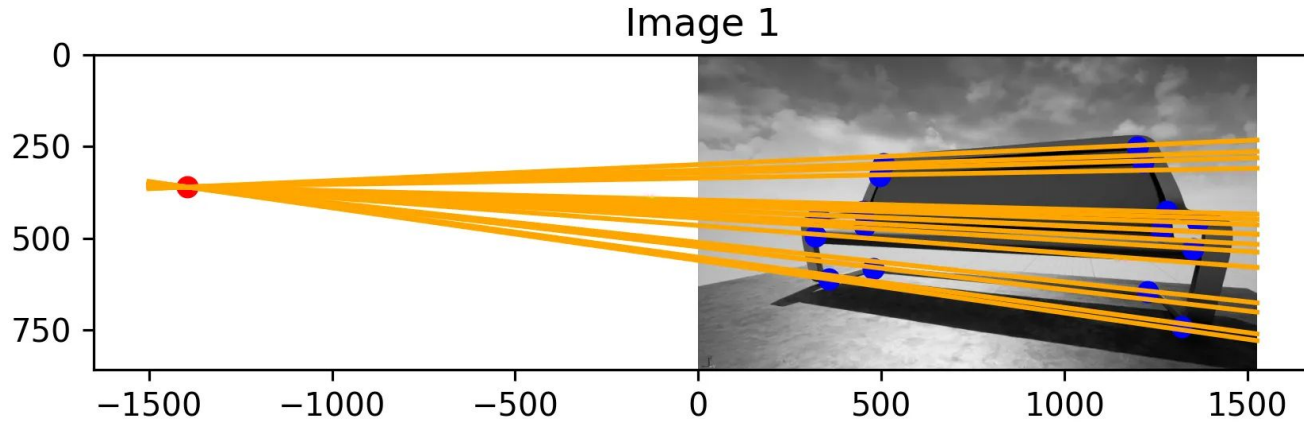
Buscar  
“features”  
(SIFT, ORB, etc.)



Calcular  
“vanishing point”

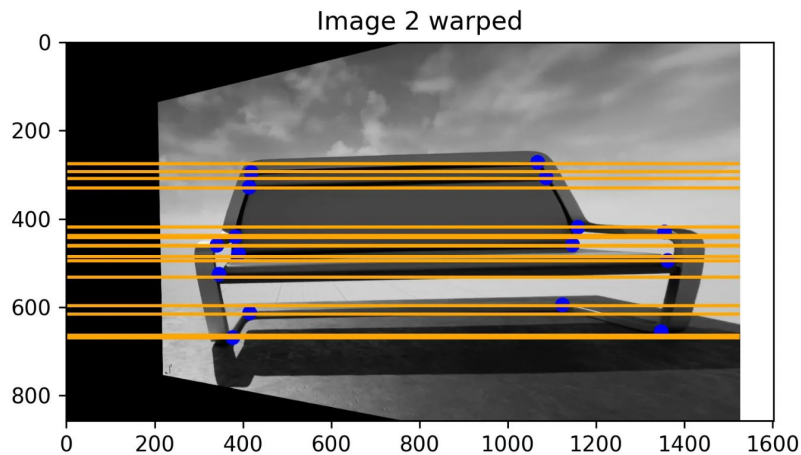
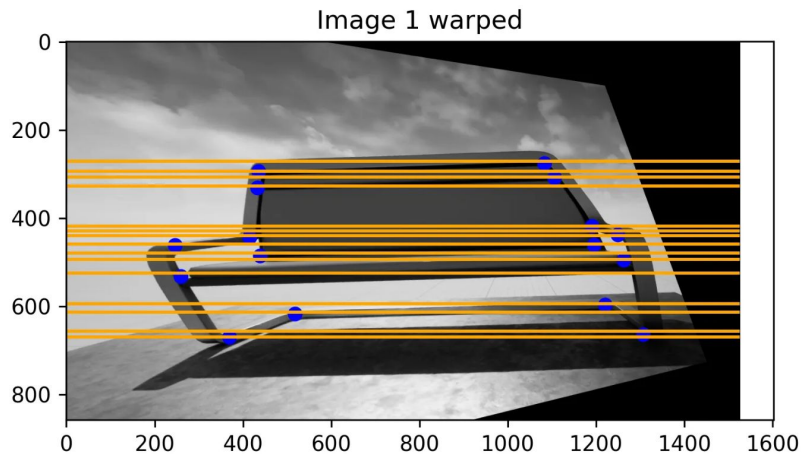
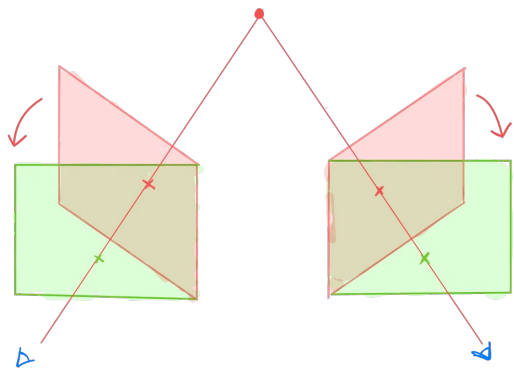


# Calculando Vanishing Points



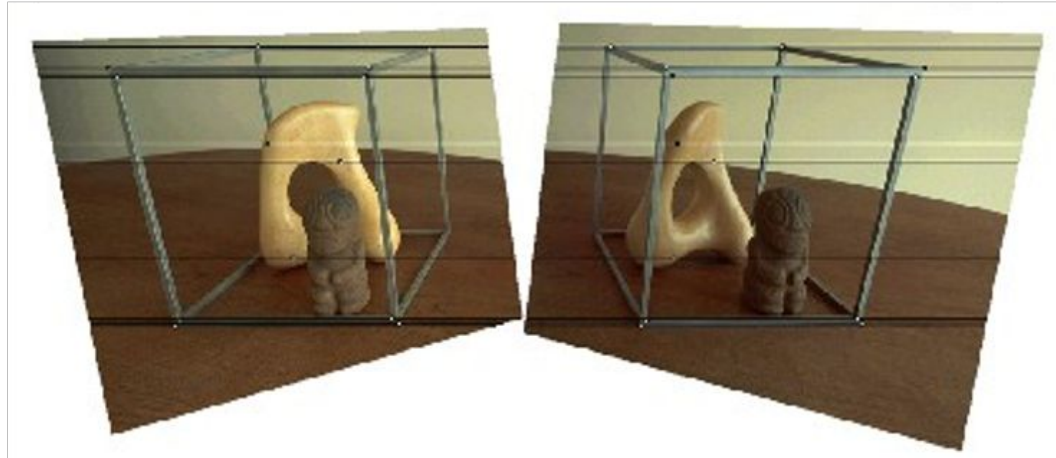
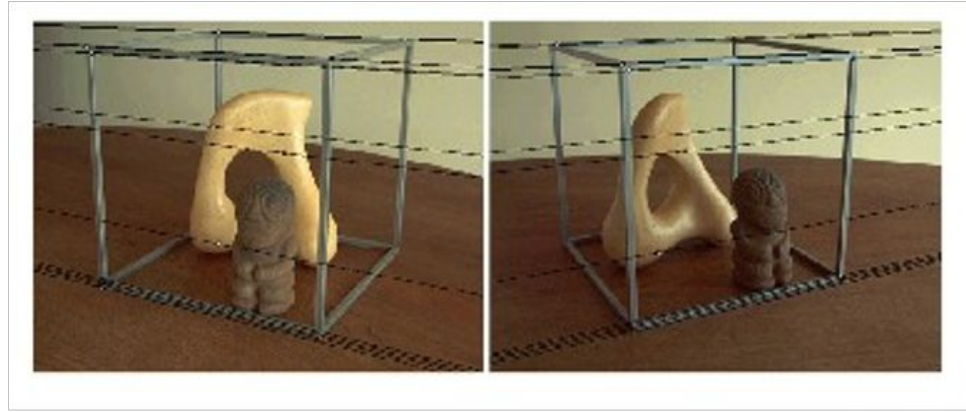
# Rectificando

- Se busca la matriz de proyección (**Matriz Esencial y Fundamental**) usando la restricción epipolar
- Se hace una transformación (Warping) de las imágenes para obtener líneas paralelas



**Nota: Con las imágenes rectificadas podemos calcular disparidad**

# Rectificando (Ejemplo 2)

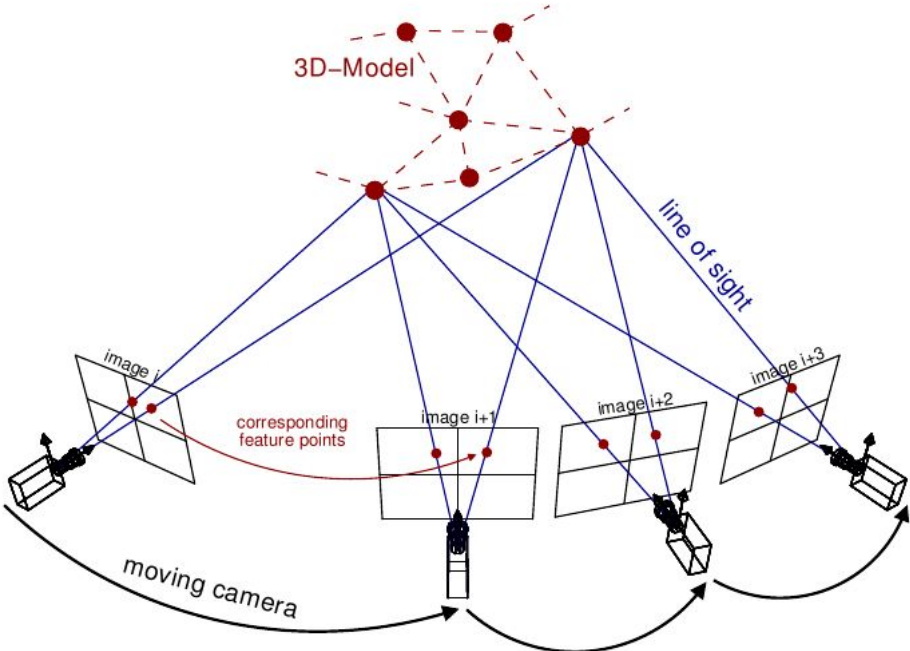
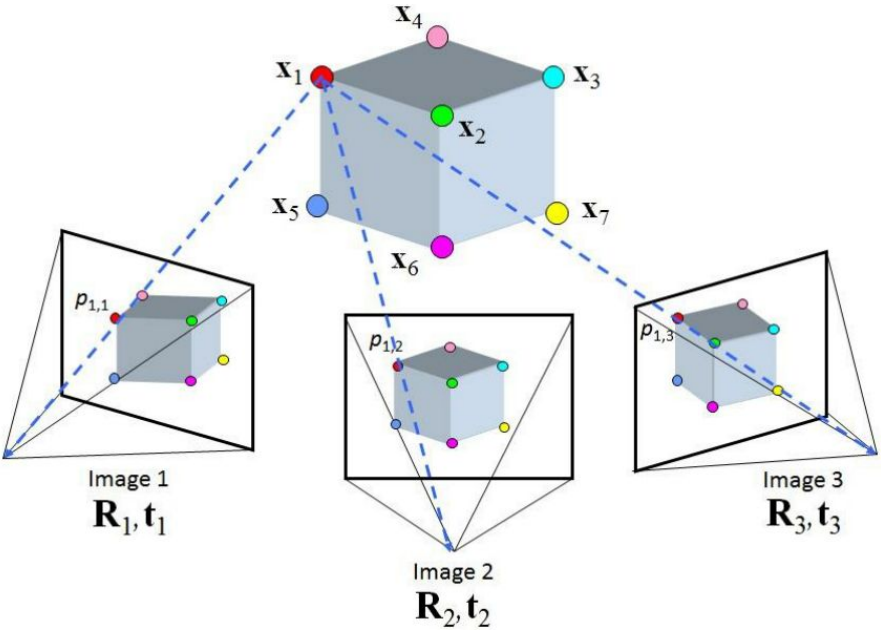


**Nota: Con las imágenes rectificadas podemos calcular disparidad**

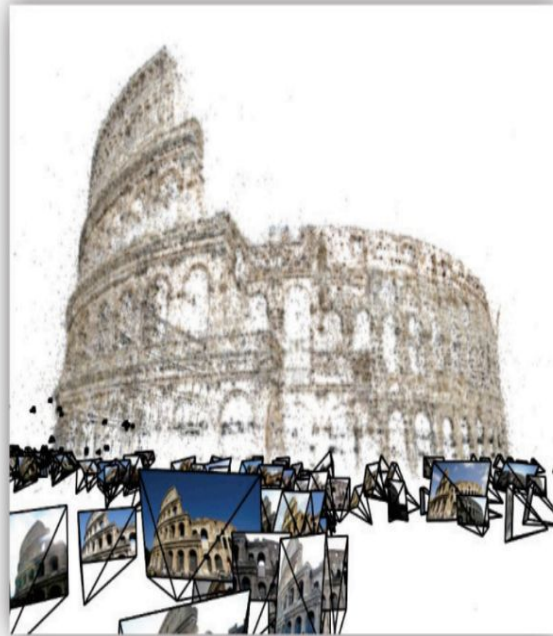


Caso más complejo  
(Structure from Motion)

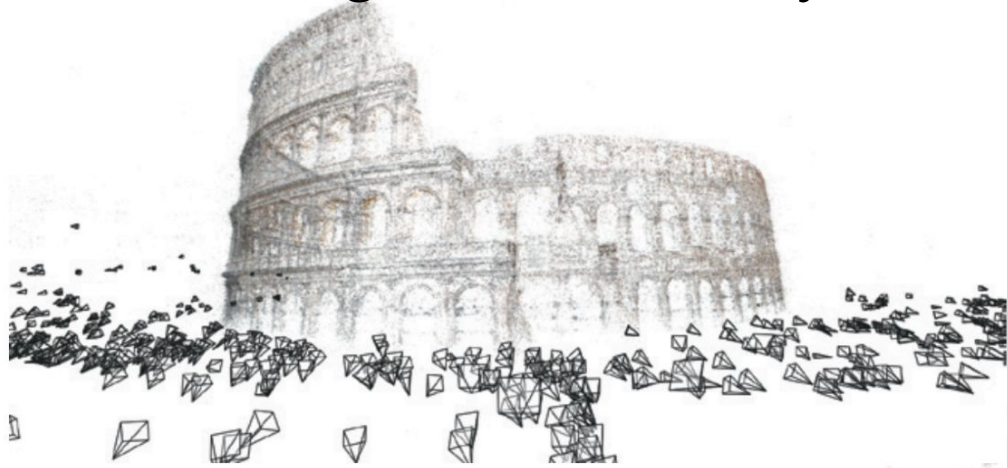
# Estructura a partir del movimiento (SfM)



# SfM: Building Rome in a Day



# SfM: Building Rome in a Day



Colosseum – 2,106 photos

Trevi Fountain – 1,936 photos



La reconstrucción de Roma (150.000 fotografías) tomó 26 horas (18 horas para hacer *matching* y 8 horas para la reconstrucción) utilizando 496 procesadores.





# Visión Estereo *in the wild*

SfM



# 4. Hands-on: Stereo Vision



# Actualicemos el repositorio!

main

1 Branch 0 Tags

Go to file

t

Add file

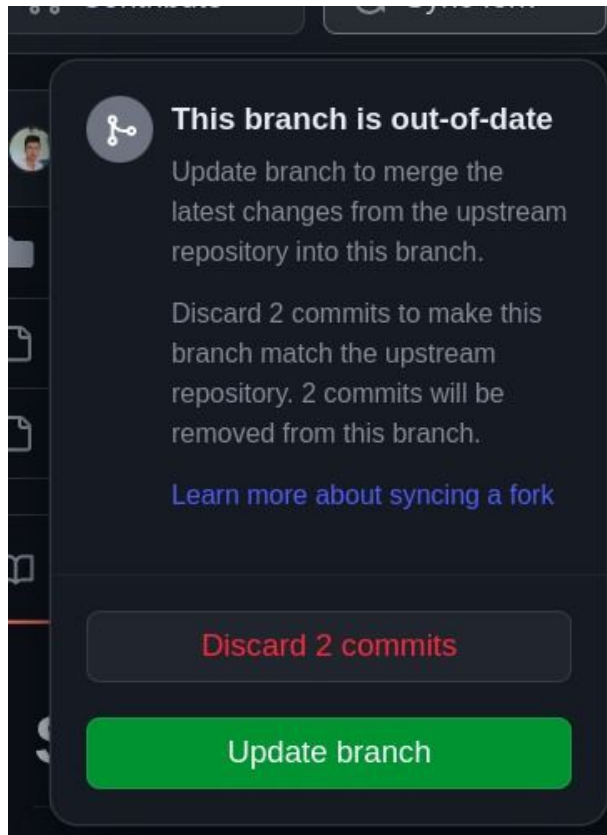
Code

This branch is 2 commits ahead of, 3 commits behind `semilleroCV/Hands-on-Computer-Vision:main`.

Contribute

Sync fork

# Actualicemos el repositorio!



**This branch is out-of-date**

Update branch to merge the latest changes from the upstream repository into this branch.

Discard 2 commits to make this branch match the upstream repository. 2 commits will be removed from this branch.

[Learn more about syncing a fork](#)

**Discard 2 commits**

**Update branch**

# Actualicemos el repositorio!

main

1 Branch 0 Tags

Go to file

t

Add file

Code

This branch is 2 commits ahead of, 3 commits behind `semilleroCV/Hands-on-Computer-Vision:main`.

Contribute

Sync fork