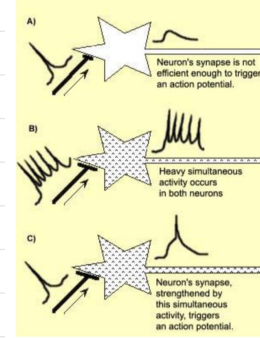


Early History

1943 McCulloch & Pitts Neuron

1949 D. Hebb: Synaptic Learning

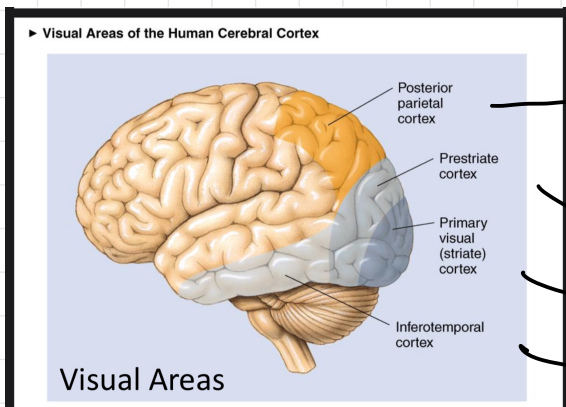


1950 Turing's suggestion: produce a programme to simulate child

~1960 McCarthy, Minsky, Newell: Represent & reason in first-order logic

Rosenblatt, Widrow: Pattern Recognition, learning

Bellman, Kalman: Estimation & Control (Markov Decision ...)



→ Spatial awareness; Movement planning;
Sensory info Integration

② V₂ → detect complex visual features (edge, contour, texture)

① V₁ → extract fundamental visual features for subsequent steps

→ visual object recognition

1962 Hubel & Wiesel: orientation sensitive neurons in V₁

1980 Fukushima: Neural Network model for pattern recognition
↳ Lack the effectiveness of backpropagation ⇒ x scale output
(unsupervised)

1989 Yan LeCun Backpropagation to train weights for handwritten digits

⊗ with GPUs, ImageNet

2012 Krizhevsky, Sutskever & Hinton ++ CV obj detection benchmark

3R: Reorganization, Recognition, Reconstruction

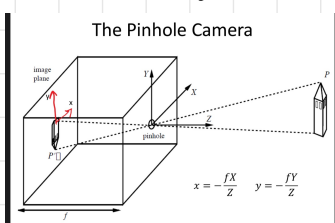
6 lessons from baby for embodied intelligence.

- multi-modal
- be incremental
- be physical
- Explore
- Be social
- language

Image Formation

Image $I(x,y)$ measures how much light is captured at pixel (x,y)

- ① where does a point (X,Y,Z) in the world get imaged
- ② what's the brightness at the resulting point (x,y)



Projection model projective transformation

① perspective projection: mapping points from [depth variation] 3D space to rays through proj center

Parallel lines converge to a vanishing point

$$\begin{bmatrix} x \\ y \\ z \end{bmatrix} = \begin{bmatrix} A_x \\ A_y \\ A_z \end{bmatrix} + \lambda \begin{bmatrix} D_x \\ D_y \\ D_z \end{bmatrix} \quad \lambda \rightarrow \infty \Rightarrow x = \frac{fX}{Z} = \frac{A_x + \lambda D_x}{\lambda D_z} = \frac{D_x}{D_z}$$

Lines perpendicular to the camera optic axis will not converge parallel to the image plane

- Farther, smaller. $1/z$
- Tilted with respect to line of sight, smaller. \cos foreshortened

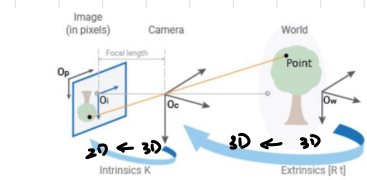
② Orthogonal projection: depth vary not much affine transformation

Camera Calibration

Pinhole camera to real cameras imaging 3D points

↓ silicon-based sensor
convert light into electrical voltage

CCD ↑ power sci
CMOS ↑ noise life



$[R|t]$ Extrinsic param: location of the camera in the 3-D scene

K Intrinsic param: optical center and the focal length of the camera

$$W \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = P \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$$

Scale factor Image points World points

Camera matrix Extrinsic Intrinsic matrix

Rotation and Translation

$[c_x \ c_y]$ - Optical center (the principal point), in pixels.

$(f_x \ f_y)$ - Focal length in pixels.

$f_x = f/p_x$

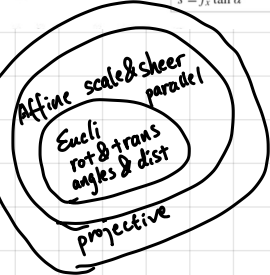
$f_y = f/p_y$

F - Focal length in world units, typically expressed in millimeters.

$(p_x \ p_y)$ - Size of the pixel in world units.

s - Skew coefficient, which is non-zero if the image axes are not perpendicular.

$s = f_x \tan \alpha$



Eucli: $\varphi(a) = Aa + t, A^T A = I$

Affine: $\varphi(a) = Aa + t, \det(A) \neq 0$

Homo: finite. infinite

Projective line: $\begin{bmatrix} x \\ 1 \end{bmatrix}, \begin{bmatrix} 2x \\ 2 \end{bmatrix}, \begin{bmatrix} b+3x \\ b+2 \end{bmatrix} \dots \begin{bmatrix} x \\ 0 \end{bmatrix}$

Projective plane: $\begin{bmatrix} \lambda x \\ \lambda y \\ \lambda \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$

$\begin{bmatrix} x \\ y \end{bmatrix} \leftrightarrow \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$ Affine

$\begin{bmatrix} x' \\ y' \\ w' \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} & t_x \\ a_{21} & a_{22} & t_y \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$

Perspective proj: $\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1/f & 0 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix} = \begin{bmatrix} x \\ y \\ z/f \end{bmatrix}$

Calibration method.

① Tsai's method: $n \geq 6$ control points 3D position + 2D coordinates in image

$$\begin{bmatrix} x_i & y_i & z_i \\ u_i & v_i \end{bmatrix}_{i \geq 6} \xrightarrow{\text{Direct Linear Transform}} P = \begin{bmatrix} p_{11} & p_{12} & p_{13} & p_{14} \\ p_{21} & p_{22} & p_{23} & p_{24} \\ p_{31} & p_{32} & p_{33} & p_{34} \end{bmatrix}_{3 \times 4} \Rightarrow \lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = P \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix}$$

Q $QP = 0$ $U, S, V = \text{svd}(Q)$ $P = V(:, 1:2)$ $K, R, T = \text{qr}(P)$

Metric: reprojection error = $\|p^i - \pi(p_{in}^i, K, R, T)\|$

Refine by minimizing K, R, T , len distort = $\argmin_{K, R, T, \text{len}} \|p^i - \pi(p_{in}^i, K, R, T)\|$

Levenberg-Marquardt

② Zhang's method: multi-views of a planar grid.

2D-2D Homography

$$\lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \begin{bmatrix} h_{11} & h_{12} & h_{13} \\ h_{21} & h_{22} & h_{23} \\ h_{31} & h_{32} & h_{33} \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} \quad QH = 0 \geq 4$$

$n \geq 4$ non-collinear

20-50 views for KRT

$H = \text{SVD}$

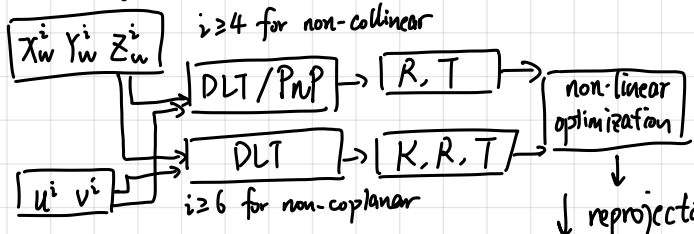
After calibration, camera localization $P_n P$

determine camera pose (R, T) relative to World Frame.

$3(P_3 P) + 1$ for disambiguation.

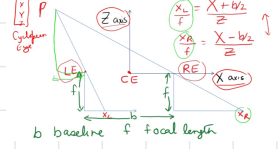
Calibrated: DLT or $P_n P$

Uncalibrated: DLT



Triangulation

Parallel Optical Axes (fixation at infinity)



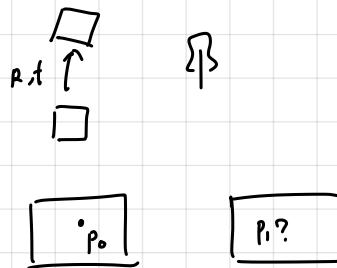
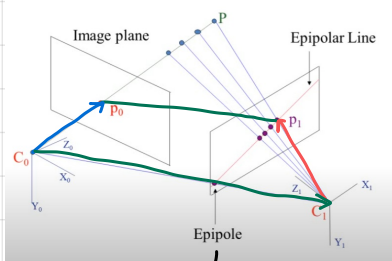
disparity = $x_L - x_R = \frac{bf}{Z}$

$Z = \frac{bf}{x_L - x_R}$

Multi-view geometry \Rightarrow reconstruct internet images

$\begin{cases} X \text{ camera} & X \text{ 3D shape} \\ \vee & N \text{ correspondences} \end{cases} \Rightarrow \text{solve for camera \& depth of pts}$

① re { camera points } → epipolar geo structure from motions



estimate the motion

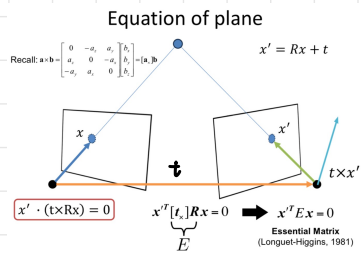
- ① in the image, forward & backward
- ② outside the image, sideways
- ③ image center, pure translation, no rot

$\vec{C_0P_0}$, $\vec{C_0C_1}$, $\vec{C_1P_1}$ are coplanar vectors

Epipolar geo: { two camera centers
two image projections of any 3D point } ⇒ on the same plane.

⇒ reduce search space

② get camera from points using epipolar geo



$$\begin{aligned} x' \cdot (t \times x) &= 0 \\ x' \cdot (t \times (Rx + t)) &= 0 \\ x' (t \times Rx + t \times t) &= 0 \\ x' (t \times Rx) &= 0 \\ x'^T \underbrace{\begin{bmatrix} t_2 & -t_1 & 0 \\ t_3 & 0 & -t_1 \\ 0 & t_3 & t_2 \end{bmatrix}}_E R x &= 0 \end{aligned}$$

$$\begin{aligned} \text{skew-symmetric orthogonal} \\ [x' \ y' \ z'] \begin{bmatrix} 0 & -t_2 & t_1 \\ t_2 & 0 & -t_1 \\ -t_1 & t_1 & 0 \end{bmatrix} \begin{bmatrix} r_{11} & r_{12} & r_{13} \\ r_{21} & r_{22} & r_{23} \\ r_{31} & r_{32} & r_{33} \end{bmatrix} \begin{bmatrix} x \\ y \\ z \end{bmatrix} &= 0 \\ [x' \ y' \ z'] \begin{bmatrix} e_{11} & e_{12} & e_{13} \\ e_{21} & e_{22} & e_{23} \\ e_{31} & e_{32} & e_{33} \end{bmatrix} \begin{bmatrix} x \\ y \\ z \end{bmatrix} &= 0 \\ \text{SVD} \\ E &\rightarrow T, R \checkmark \end{aligned}$$