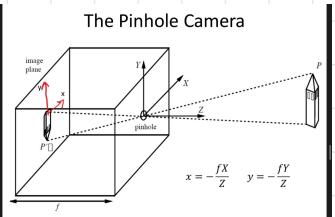


Image Formation

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

$$\begin{bmatrix} X \\ Y \\ Z \end{bmatrix} = \begin{bmatrix} Ax \\ Ay \\ Az \end{bmatrix} + \lambda \begin{bmatrix} Dx \\ Dy \\ Dz \end{bmatrix} \quad \lambda \rightarrow \infty \Rightarrow X = \frac{fX}{Z} = \frac{Ax + \lambda Dx}{Az} = \frac{Dx}{Dz}$$

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

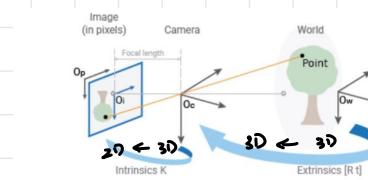
- Farther, smaller. $1/z$
- Tilted with respect to line of sight, smaller. \cos Foreshortened

- ② Orthogonal projection: depth vary not much **affine transformation**

Camera Calibration

Pinhole camera to real cameras imaging 3D points

↓ silicon-based sensor ↑ power sci
convert light into electrical voltage ↑ CMOS ↑ noise life

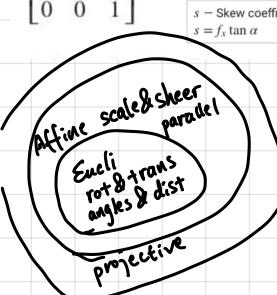


[at] Extrinsic param: location of the camera in the 3-D scene

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

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

$[c_x \ c_y]$ – 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.
$s = f_x \tan \alpha$



$$a_1 X_1 + a_2 X_2 + a_3 X_3 = 0$$

$$\text{point on line: } \begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix} \begin{bmatrix} X_1 \\ X_2 \\ X_3 \end{bmatrix} = 0$$

$$\text{intersect: } \begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix} \times \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix} = \begin{bmatrix} x \\ y \\ 0 \end{bmatrix}$$

$$\Rightarrow \text{parallel}$$

$$\begin{bmatrix} x \\ y \\ 0 \end{bmatrix} \xleftarrow{\text{Affine}} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$\begin{bmatrix} x \\ y \\ 0 \end{bmatrix} \xleftarrow{\text{Perspective Proj}} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1/f \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = \begin{bmatrix} x \\ y \\ z/f \end{bmatrix}$$

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

$$\begin{bmatrix} X_i & Y_i & Z_i \\ U_i & V_i \end{bmatrix}_{i=1}^6 \xrightarrow{\text{Direct Linear Transform}} P: \begin{bmatrix} P_{11} & P_{12} & P_{13} & P_{14} \\ P_{21} & \dots & \dots & P_{24} \\ \vdots & \vdots & \vdots & \vdots \end{bmatrix}_{3 \times 4} \Rightarrow \lambda \begin{bmatrix} U \\ V \\ 1 \end{bmatrix} = P \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$$

$$Q \quad U, S, V = \text{svd}(Q) \quad P = V C^{-1} U^T \quad K, R, T = \text{genc} P$$

$$\text{Metric: reprojection error} = \|P^i - \pi(P_{\text{inv}}, K, R, T)\|$$

$$\text{Refine by minimizing } K, R, T, \text{ len distor} = \underset{K, R, T, \text{len}}{\text{argmin}} \|P^i - \pi(P_{\text{inv}}, K, R, T)\|$$

Levenberg - Marquardt

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

$Z=0$ reduce coupling

2D-2D Homography

$$\lambda \begin{bmatrix} U \\ V \\ 1 \end{bmatrix} = \begin{bmatrix} h_{11} & h_{12} & h_{13} \\ h_{21} & h_{22} & h_{23} \\ h_{31} & h_{32} & h_{33} \end{bmatrix} \begin{bmatrix} X \\ Y \\ 1 \end{bmatrix}$$

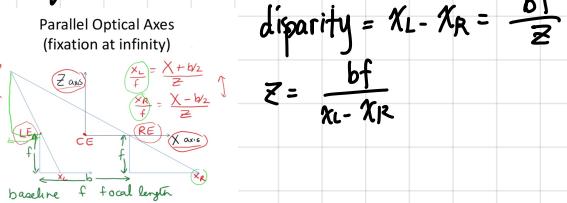
20-50 views for KRT

After calibration, camera localization PnP

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

3 (P3P) + 1 for disambiguation.

Triangulation



$$\text{disparity} = X_L - X_R = \frac{bf}{z}$$

$$z = \frac{bf}{X_L - X_R}$$

Multi-view geometry → reconstruct internet images

3 camera 3D shape \Rightarrow solve for camera & depth of pts
N correspondences

$$\text{reflection} \begin{bmatrix} \cos 2\theta & \sin 2\theta \\ \sin 2\theta & -\cos 2\theta \end{bmatrix}$$

Problem: Explain the concepts of vanishing points and the vanishing line for a plane in 3D under perspective projection.

Solution: Plane Equation: Consider a plane in 3D given by

$$N_x X + N_y Y + N_z Z = d,$$

where $N = (N_x, N_y, N_z)$ is the unit normal, and (X, Y, Z) are world coordinates on the plane. Lines on the Plane: A line in the plane is parametrized as

$$X(t) = X_0 + tD,$$

where X_0 is a point on the plane, and D is the direction vector, satisfying $N \cdot D = 0$.

Projection and Vanishing Point: Under perspective projection, a vanishing point arises as

$$\lim_{t \rightarrow \infty} (\text{image of } X_0 + tD),$$

corresponding to a point at infinity in direction D in homogeneous coordinates.

Vanishing Line Condition: All directions D in the plane satisfy $N \cdot D = 0$. Their vanishing points form the vanishing line, expressed in homogeneous coordinates as

$$N_x x + N_y y + N_z z = 0.$$

Conclusion: Every line on the plane has a vanishing point on this line, defining the vanishing associated with the plane.

Algorithm 4 Class-Conditional Sampling

```

1: input: one-hot vector c, classifier guidance scale γ
2:  $x_t = x_0 \sim \mathcal{N}(0, I)$ 
3: for t from 0 to 1, step size  $\frac{1}{T}$  do
4:    $u_{\text{uncond}} = u_\theta(x_t, t, 0)$ 
5:    $u_{\text{cond}} = u_\theta(x_t, t, c)$ 
6:    $u = u_{\text{uncond}} + \gamma(u_{\text{cond}} - u_{\text{uncond}})$ 
7:    $x_t = x_t + \frac{1}{T}u$ 
8: end for
9: return  $x_t$ 

```

Center the points: Define

$$u'_j = u_j - \bar{u}, \quad v'_j = v_j - \bar{v}.$$

Compute the cross-covariance matrix:

$$H = \sum_{j=1}^4 u'_j (v'_j)^T = \begin{pmatrix} 0 & -9 \\ 2 & 0 \end{pmatrix}.$$

Obtain the rotation: Using the SVD of $H = U\Sigma V^T$, the optimal rotation is

$$R = VU^T = \begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix},$$

a 90° rotation. Compute the translation:

$$t = \bar{v} - R\bar{u} = \begin{pmatrix} 0 & \frac{3}{4} \end{pmatrix}.$$

Thus, the best Euclidean transformation is

$$E(u) = \begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix} u + \begin{pmatrix} 0 \\ 0.75 \end{pmatrix}.$$

① Epipolar geometry: $\vec{x}, \vec{x}', \vec{t}$ coplanar $\vec{x} \cdot (\vec{t} \times \vec{x}) = 0$

② Get camera: $\vec{x}_{\text{pixel}} \vec{K}^{-1} \vec{T}^T \vec{R} \vec{K}^{-1} \vec{x}_{\text{pixel}} = 0$

$$F \xrightarrow{\vec{K}^T \vec{F} \vec{K}'} \xrightarrow{\text{SVD}} T, R$$

③ Triangulate for depth $\begin{cases} \vec{x}_{\text{pixel}} = K\vec{x} \\ \text{Linear Least Square} \end{cases}$

$\vec{x}_{\text{pixel}} = K'\vec{x}' = K'(R\vec{x} + \vec{T})$

skew-symmetric orthogonal

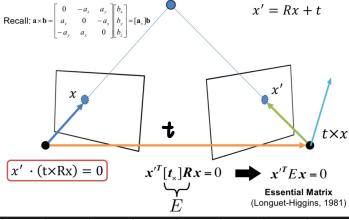
$$\begin{array}{l} \left[\begin{array}{c} \vec{x} \\ \vec{y} \\ \vec{z} \end{array} \right] \left[\begin{array}{ccc} 0 & -tx & ty \\ tx & 0 & -tx \\ -ty & tx & 0 \end{array} \right] \left[\begin{array}{ccc} f_{11} & f_{12} & f_{13} \\ f_{21} & f_{22} & f_{23} \\ f_{31} & f_{32} & f_{33} \end{array} \right] \left[\begin{array}{c} \vec{x} \\ \vec{y} \\ \vec{z} \end{array} \right] = 0 \quad \text{bundle adj} \\ \downarrow \\ \left[\begin{array}{c} \vec{x} \\ \vec{y} \\ \vec{z} \end{array} \right] \left[\begin{array}{ccc} 0 & -tx & ty \\ tx & 0 & -tx \\ -ty & tx & 0 \end{array} \right] \left[\begin{array}{ccc} f_{11} & f_{12} & f_{13} \\ f_{21} & f_{22} & f_{23} \\ f_{31} & f_{32} & f_{33} \end{array} \right] \left[\begin{array}{c} \vec{x} \\ \vec{y} \\ \vec{z} \end{array} \right] = 0 \quad \text{sparse reconstruct} \\ \downarrow \\ \vec{x}_{\text{pixel}}^T \vec{K}^{-1} E \vec{K}^{-1} \vec{x}_{\text{pixel}} = 0 \quad \text{depth estim} \quad \text{MVS} \quad \text{dense reconstruct} \\ \downarrow \\ \text{merge depth map} \quad \text{mesh} \end{array}$$

epipolar line in 2nd image given \vec{x}

$$\ell' = F\vec{x} \quad \ell = F^T\vec{x}'$$

\Rightarrow reduce search space

Equation of plane



Epipole estimate the motion

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

$$\begin{bmatrix} u \\ v \end{bmatrix} = \begin{bmatrix} \dot{x} \\ \dot{y} \end{bmatrix} = \frac{1}{2} \begin{bmatrix} -1 & 0 & x \\ 0 & -1 & y \end{bmatrix} \begin{bmatrix} tx \\ ty \\ tz \end{bmatrix} + \begin{bmatrix} xy - (1+x^2) & -y \\ 1+y^2 & -xy \end{bmatrix} \begin{bmatrix} ux \\ vy \\ wz \end{bmatrix}$$

Assume that the camera moves with translational velocity $t = (t_x, t_y, t_z)$ and angular velocity $\omega = (\omega_x, \omega_y, \omega_z)$. Eq.(3.1) is used to characterize the movement of X,

$$\dot{X} = -t - \omega \wedge X, \quad (3.1)$$

which can be written out in coordinates as Eq.(3.2):

$$\begin{bmatrix} \dot{X} \\ \dot{Y} \\ \dot{Z} \end{bmatrix} = \begin{bmatrix} t_x \\ t_y \\ t_z \end{bmatrix} - \begin{bmatrix} \omega_y Z - \omega_z Y \\ \omega_z X - \omega_x Z \\ \omega_x Y - \omega_y X \end{bmatrix}. \quad (3.2)$$

$$\begin{aligned} \dot{x} &= \frac{\dot{X}Z - \dot{Z}X}{Z^2} \quad y = \frac{\dot{Y}Z - \dot{Z}Y}{Z^2} \quad \text{as } x = \frac{X}{Z} \quad y = \frac{Y}{Z} \\ &= \frac{[-tx - w_y z + w_z y]}{Z^2} Z - \frac{[-t_z - w_x z + w_y x]}{Z^2} Z \\ &= \frac{1}{Z} (-tx + xt_z) + xyw_x - (1+x^2)wy + yw_z. \end{aligned}$$

Light field

Amount of light flowing in every direction $L(x, y, z, \theta, \phi) \Rightarrow (C, \sigma)$

$$\text{power} : L(p_i, \vec{p}_i) dA \cos \theta \quad S = \theta r$$

Lambertian model: easy version, color x change based on view direction

$$L = \rho \lambda (n \cdot s) \quad \begin{cases} \rho = 0 & \text{absorb} \\ \rho = 1 & \text{reflect} \end{cases}$$

cones / rods are pixels in eye
response: weak light at peak sensitivity
 \Leftrightarrow strong light at low sensitivity

$$= \int L(\lambda) R(\lambda) d\lambda$$

brightness across the edge changes a lot.
reflectance / illumination / surface geo(depth)

$$\Rightarrow \text{Bidirectional Reflectance Distribution Function.}$$

$$f(w_i, w_r) = \frac{dL_r(w_r)}{dE_i(w_i)} = \frac{\text{radiance}}{\text{irradiance}}$$

good design of visual systems \Rightarrow hierarchical & feedforward

Problem: Verify Rodrigues' formula by considering the powers of the skew-symmetric matrix associated with the cross product with a vector.

Solution: Let $s = (s_1, s_2, s_3)$ be a unit vector and define its associated skew-symmetric matrix

$$\hat{s} = \begin{pmatrix} 0 & -s_3 & s_2 \\ s_3 & 0 & -s_1 \\ -s_2 & s_1 & 0 \end{pmatrix}.$$

One can show that

$$\hat{s}^2 = ss^T - I, \quad \text{and} \quad \hat{s}^3 = -\hat{s}.$$

Expanding the matrix exponential in a Taylor series,

$$e^{\phi \hat{s}} = I + \phi \hat{s} + \frac{\phi^2}{2!} \hat{s}^2 + \frac{\phi^3}{3!} \hat{s}^3 + \dots,$$

and grouping even and odd terms using the identities above, one obtains

$$e^{\phi \hat{s}} = I + \sin \phi \hat{s} + (1 - \cos \phi) \hat{s}^2.$$

This is Rodrigues' formula for constructing the rotation matrix.

Filtering | Gaussian Filter smooth x preserve edge / convolve an image
Anisotropic diffusion edge A comp expensive / High dif in uni + low dif on edge
Bilateral Filter edge & comp exp for t kernel / ↑ weights to similar pixels
Non-Local Means keep details & slow / avg similar patches across the image

\Rightarrow data driven kernel + Full field of view

Algorithm 1 Training def homography_solve(u, v):

```
N = u.shape[1]
A = []
b = []

for i in range(N):
    x, y = u[0, i], u[1, i]
    x_p, y_p = v[0, i], v[1, i]

    A.append([x, y, 1, 0, 0, 0, -x_p*x, -x_p*y])
    b.append(x_p)

A.append([x, y, 1, 0, 0, 0, -y_p*x, -y_p*y])
b.append(y_p)
```

A = np.array(A, dtype=np.float64)
b = np.array(b, dtype=np.float64)

Solve
ATA = A.T @ A
ATb = A.T @ b

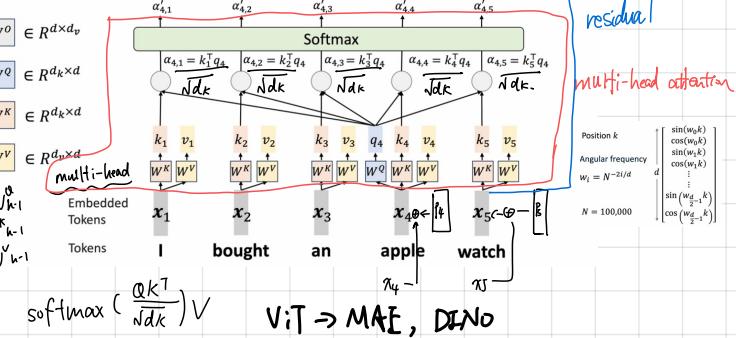
Algorithm 3 Class-Conditional Training

```
1: repeat
2:   x1, c ~ clean image and label from training set
3:   Make c into a one-hot vector
4:   with probability pmodel set c to zero-vector.
5:   t ~ Uniform([0, 1])
6:   x0 ~ N(0, I)
7:   x_t = (1 - t)x0 + tx1
8:   Take gradient descent step on
   ∇θ ||(x1 - x0) - uθ(x_t, t)||2
9: until happy
```

return H

$$\text{Updated feature } x'_4 = W^0 (a'_{4,1} v_1 + a'_{4,2} v_2 + a'_{4,3} v_3 + a'_{4,4} v_4 + a'_{4,5} v_5)$$

\Rightarrow feed forward, layer norm



$$\text{softmax} \left(\frac{QK^T}{Ndk} \right) V$$

ViT \Rightarrow MAE, DINO

Variational Autoencoders

- The Encoder maps x to z .
- But need to be able to sample from z like $p(z)$
- So we need the encoder to learn $p(z|x)$
- But how to get $p(z|x)$...?
- What's wrong?
- So learn $q(z|x)$ and make it close to $p(z)$ $DL(q(z|x)||p(z))$
- Through ELBO (will do later):

$$\max \log p(x) = \mathbb{E}_{q(z|x)} \log p(x|z) - DL(q(z|x)||p(z))$$

i.e. max $p(x|z)$ while regularizing $q(z|x)$ to be close to $p(z)$

Objective: $\mathbb{E}_{z \sim q(z|x)} \log p(x|z) - DL(q(z|x)||p(z))$

$$D_{KL}(q_\phi(z|x) || p(z)) = \mathbb{E}_{z \sim q_\phi(z|x)} \left[\log \frac{q_\phi(z|x)}{p(z)} \right]$$

$$= \frac{\mathbb{E}_{z \sim q_\phi(z|x)} [\log q_\phi(z|x)] - \mathbb{E}_{z \sim p(z)} [\log p(z)]}{H(q_\phi(z|x))}$$

$$H(q_\phi(z|x)) = - \mathbb{E}_{z \sim p(z)} [\log p(z)] = - \int p(z) \log p(z) dz$$

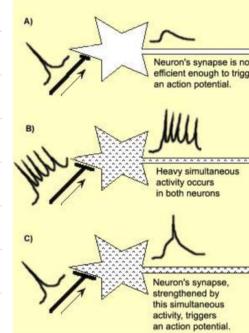
$$- D_{KL}(q_\phi(z|x) || p(z)) = \mathbb{E}_{z \sim q_\phi(z|x)} [\log p(z)] + H(q_\phi(z|x))$$

Re-written Objective: $\mathbb{E}_{z \sim q(z|x)} \log p(x|z) + \log p(z) + H(q(z|x))$

Early History

1943 McCulloch & Pitts Neuron

1949 D. Hebb: Synaptic Learning

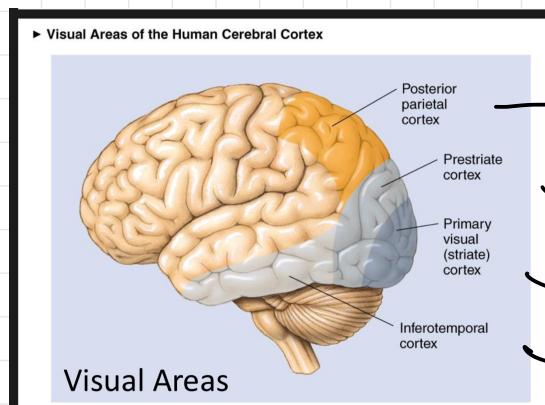


1950 Turing's suggestions produce a programme to simulate child

~1960 McCarthy, Minsky, Newell: Rensselaer & 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 \Rightarrow x scale output
(Unsupervised)
abstract repres.-
i.e. learned feature
activations

1989 Yann LeCun Backpropagation to train weights for handwritten digits
↗ with GPUs, Imagenet

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

3R : Reorganization, Recognition, Reconstruction

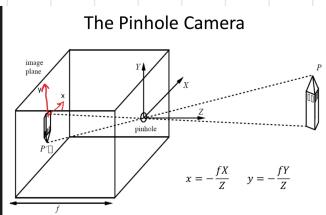
6 lessons from baby for embodied intelligence.

- 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

$$\begin{bmatrix} X \\ Y \\ Z \end{bmatrix} = \begin{bmatrix} Ax \\ Ay \\ Az \end{bmatrix} + \lambda \begin{bmatrix} Dx \\ Dy \\ Dz \end{bmatrix} \quad \lambda \rightarrow \infty \Rightarrow X = \frac{fX}{Z} = \frac{Ax + \lambda Dx}{Az} = \frac{Dx}{Dz}$$

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

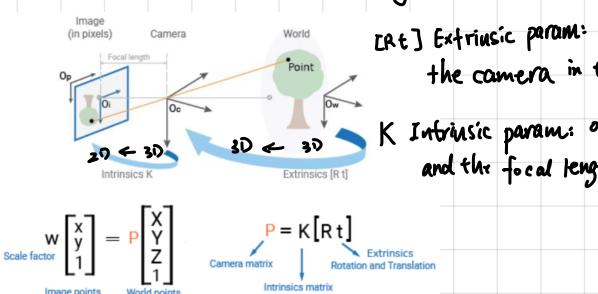
- Farther, smaller. $1/Z$
- Tilted with respect to line of sight, smaller. \cos Foreshortened

- ② Orthogonal projection: depth vary not much affine transformation

Camera Calibration

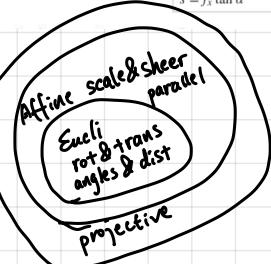
Pinhole camera to real cameras imaging 3D points

↓ silicon-based sensor ↑ power sci
convert light into electrical voltage ↑ CMOS ↑ noise life



$W \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = P \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$	$P = K[R t]$
Scale factor	Intrinsics matrix

$[x, y]$ – 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.
 $s = f_x \tan \alpha$



$$a_1 X_1 + a_2 X_2 + a_3 X_3 = 0$$

$$\text{point on line: } \begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = 0$$

$$\text{intersect: } \begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix} \times \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix} = \begin{bmatrix} x \\ y \\ 0 \end{bmatrix} \rightarrow \text{parallel}$$

$$\text{Affine: } \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} \leftrightarrow \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$\text{Perspective Proj: } \begin{bmatrix} x' \\ y' \\ w \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} & t_x \\ a_{21} & a_{22} & t_y \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ z/f \end{bmatrix} = \begin{bmatrix} x \\ y \\ z/f \end{bmatrix}$$

Calibration method.

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

$$\begin{bmatrix} X_i & Y_i & Z_i \\ U_i & V_i \end{bmatrix}_{i \geq 6} \xrightarrow{\text{Direct Linear Transform}} P = \begin{bmatrix} P_{11} & P_{12} & P_{13} & P_{14} \\ P_{21} & P_{22} & P_{23} & P_{24} \\ \vdots & \vdots & \vdots & \vdots \\ P_{n1} & P_{n2} & P_{n3} & P_{n4} \end{bmatrix}_{3 \times 4} \Rightarrow \lambda \begin{bmatrix} U \\ V \\ 1 \end{bmatrix} = P \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$$

$$Q \quad U, S, V = \text{svd}(Q)$$

$$P = V(:, 1:2)$$

$$K, R, T = \text{qr}(P)$$

Metric: reprojection error = $\|p^i - \pi(P^i, K, R, T)\|$

Refine by minimizing $K, R, T, \text{len distor} = \arg\min_{K, R, T, \text{len}} \|p^i - \pi(p^i, K, R, T)\|^2$

Levenberg-Marguardt

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

$Z = 0$ reduce complexity

Homography

$$\lambda \begin{bmatrix} U \\ V \\ 1 \end{bmatrix} = \begin{bmatrix} h_{11} & h_{12} & h_{13} \\ h_{21} & h_{22} & h_{23} \\ h_{31} & h_{32} & h_{33} \end{bmatrix} \begin{bmatrix} X \\ Y \\ 1 \end{bmatrix} \quad QH = 0 \geq 4$$

20-50 views for K, R, T

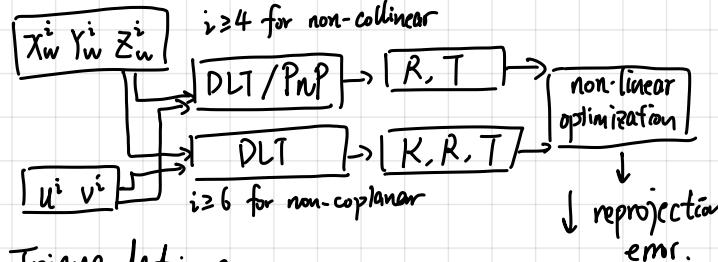
$H = \text{SVD}$ non-collinear

After calibration, camera localisation PnP

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

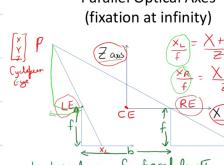
$3(P_3 P) + 1$ for disambiguation.

Calibrated = DLT or PnP Uncalibrated: DLT



Triangulation

Parallel Optical Axes (fixation at infinity)



$$\text{disparity} = X_L - X_R = \frac{bf}{Z}$$

$$Z = \frac{bf}{X_L - X_R}$$

Multi-view geometry → reconstruct internet images

$\{X \text{ camera} \quad X \text{ 3D shape} \quad \Rightarrow \text{solve for camera \& depth of pts}$
 $\checkmark N \text{ correspondences}$

$$\text{rot} \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix} \text{ reflection} \begin{bmatrix} \cos 2\theta & \sin 2\theta \\ \sin 2\theta & -\cos 2\theta \end{bmatrix}$$

① Epipolar geo: { two camera centers
two image projections of any 3D point } \Rightarrow on the same plane. Multi-view geometry

② get camera from points using epipolar geo

Equation of plane

Recall: $a \cdot b = \begin{bmatrix} 0 & -a_3 & a_2 \\ a_3 & 0 & -a_1 \\ -a_2 & a_1 & 0 \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix} \cdot \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix} = [a, b] \cdot [b]$

$$x' \cdot (t x x') = 0$$

$$x' \cdot (t x (Rx + t)) = 0$$

$$x' \cdot (txRx + tx^2) = 0$$

$$x' \cdot (txRx) = 0$$

$$x'^T \underbrace{\begin{bmatrix} t \\ Rx \end{bmatrix}}_E R x = 0$$

Essential Matrix (Longuet-Higgins, 1981)

Epipole estimate the motion

① in the image, forward & backward

② outside the image, sideways

③ image center, pure translation, no rot

$$x' \cdot (t x x') = 0 \quad [x' \ y' \ z'] \begin{bmatrix} 0 & -t_2 & t_1 \\ t_2 & 0 & -t_3 \\ -t_1 & t_3 & 0 \end{bmatrix} \begin{bmatrix} r_{11} & r_{12} & r_{13} \\ r_{21} & r_{22} & r_{23} \\ r_{31} & r_{32} & r_{33} \end{bmatrix} \begin{bmatrix} x \\ y \\ z \end{bmatrix} = 0$$

skew-symmetric orthogonal

$$x' \cdot (txRx) = 0 \quad [x' \ y' \ z'] \underbrace{\begin{bmatrix} e_{11} & e_{12} & e_{13} \\ e_{21} & e_{22} & e_{23} \\ e_{31} & e_{32} & e_{33} \end{bmatrix}}_{\text{SVD}} \begin{bmatrix} x \\ y \\ z \end{bmatrix} = 0$$

$$x'^T \underbrace{\begin{bmatrix} t \\ Rx \end{bmatrix}}_E R x = 0$$

$E \rightarrow T, R \checkmark$

Λ we don't have \vec{x}, \vec{x}' in 3D
we have corresponding 2D points

\vec{x}, \vec{x}' is in the camera coord

i.e. $x_{\text{pixel}} = K \vec{x} = K[R|T] Pw$

$\vec{x} = K^{-1} x_{\text{pixel}} \quad \vec{x}' = K^{-1} x'_{\text{pixel}}$

$\Rightarrow x'^T \underbrace{K'^T E K^{-1} x_{\text{pixel}}}_F = 0$

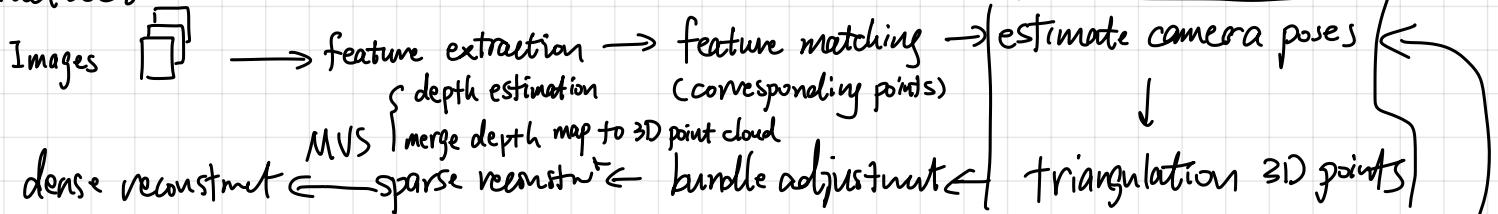
$F \xrightarrow{\text{SVD}} E \xrightarrow{\text{SVD}} T, R$

epipolar line in 2nd image given \vec{x}

$l' = Fx \quad l = F^T x'$

\Rightarrow reduce search space

Practices:



① Epipolar geometry: $\vec{x}, \vec{x}', \vec{t}$ coplanar $\vec{x}' \cdot (t x x) = 0$

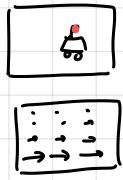
② Get camera: $x_{\text{pixel}} \underbrace{K'^T \underbrace{T \ R \ K^{-1}}_F x_{\text{pixel}} = 0}$

$F \xrightarrow{\text{SVD}} E \xrightarrow{\text{SVD}} T, R$

③ Triangulate for depth $\left\{ \begin{array}{l} x_{\text{pixel}} = K \vec{x} \\ x_{\text{pixel}} = K' \vec{x}' = K'(R \vec{x} + \vec{t}) \end{array} \right.$
Linear Least Square x best :: Noise

Dynamic Perspective

Moving camera \rightarrow scene depth



$$\text{optical flow } \begin{cases} u = \frac{\partial x}{\partial t} \\ v = \frac{\partial y}{\partial t} \end{cases}$$

$$\begin{bmatrix} u \\ v \end{bmatrix} = \begin{bmatrix} \dot{x} \\ \dot{y} \end{bmatrix} = \frac{1}{z} \begin{bmatrix} -1 & x \\ 0 & -1 \end{bmatrix} \begin{bmatrix} tx \\ ty \\ tz \end{bmatrix} + \begin{bmatrix} xy & -(1+x^2) & y \\ 1+y^2 & -xy & -x \end{bmatrix} \begin{bmatrix} w_x \\ w_y \\ w_z \end{bmatrix}$$

Assume that the camera moves with translational velocity $t = (t_x, t_y, t_z)$ and angular velocity $\omega = (\omega_x, \omega_y, \omega_z)$. Eq.(3.1) is used to characterize the movement of X ,

$$\dot{X} = -t - \omega \wedge X, \quad (3.1)$$

which can be written out in coordinates as Eq.(3.2):

$$\begin{bmatrix} \dot{X} \\ \dot{Y} \\ \dot{Z} \end{bmatrix} = - \begin{bmatrix} t_x \\ t_y \\ t_z \end{bmatrix} - \begin{bmatrix} \omega_y Z - \omega_z Y \\ \omega_z X - \omega_x Z \\ \omega_x Y - \omega_y X \end{bmatrix}. \quad (3.2)$$

$$\begin{aligned} \dot{X} &= \frac{\dot{x}z - \dot{x}z}{z^2} \quad \dot{Y} = \frac{\dot{y}z - \dot{y}z}{z^2} \quad \text{as } x = \frac{X}{z} \quad y = \frac{Y}{z} \\ &= \frac{\{-tx - w_y z + w_z y\}z - \{-t_z - w_x y + w_y x\}x}{z^2} \\ &= \frac{1}{z}(-tx + xt_z) + xyw_x - (1+x^2)wy + yw_z. \end{aligned}$$

Radiometry of Image Formation

Image is an array of brightness values. (3 arrays for RGB Images)

Pinhole camera models where a point is projected

Brightness at a pixel (x, y) ?

\Rightarrow Irradiance W/m^2 response at a pixel $\Rightarrow E \sigma A dt$

E (all directions) ↓ light intensity ↓ exposure time.

Radiance $L = \frac{E}{dA \cos \theta d\Omega}$ (in a specific direction)
the power of light travels

Light field

Amount of light flowing in every direction $L(x, y, z, \theta, \phi) \Rightarrow (C, \sigma)$

power : $L(p_i, \vec{p}_i) dA \cos \theta \frac{dA_2 \cos \theta_2}{r^2} \quad S = \theta \sigma$

Lambertian model: easy version, color x change based on view direction

$L = \rho \lambda (n, s) \quad \begin{cases} \rho = 0 & \text{absorb} \\ \rho = 1 & \text{reflect} \end{cases} \Rightarrow$ Bidirectional Reflectance Distribution Function.

$$f(w_i, w_r) = \frac{dL_r(w_r) \text{radiance}}{dE_i(w_i) \text{irradiance}} = \frac{dL_r(w_r)}{Li(w_i) \cos \theta_i d\omega_i}$$

Color Vision

cones / rods are pixels in eye

response: weak light at the peak sensitivity

\Leftrightarrow strong light at lower sensitivity

$$= \int L(\lambda) R(\lambda) d\lambda$$

power sensitivity

brightness across the edge changes a lot.

reflectance / illumination / surface geo (depth)

Human Vision

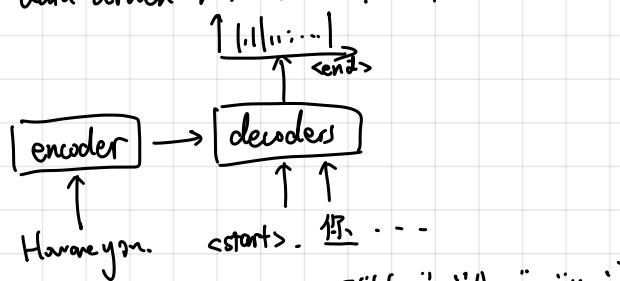
weighting function $[-1 \rightarrow 1] \Rightarrow$ receptive field

good design of visual systems \Rightarrow hierarchical & feed-forward

Sequence Model.

Filtering | Gaussian Filter smooth & preserve edge / convolve an image
| Anisotropic diffusion edge \wedge comp expensive / High dif in uni + low dif on edge
| Bilateral Filter edge \wedge comp exp for + kernel / \uparrow weights to similar pixels
| Non-Local Means keep details \wedge slow / avg similar patches across the image

\Rightarrow data driven kernel + Full field of view



① Tokenization: separate ["How", "Are", "You"]

② One-hot: \wedge X capture relation / semantic meaning

$$\begin{matrix} \text{How} \\ \# \text{token} \end{matrix} \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \quad \begin{matrix} \text{Are} \\ \# \text{token} \end{matrix} \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \quad \begin{matrix} \text{You} \\ \# \text{token} \end{matrix} \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$$

③ Token embedding

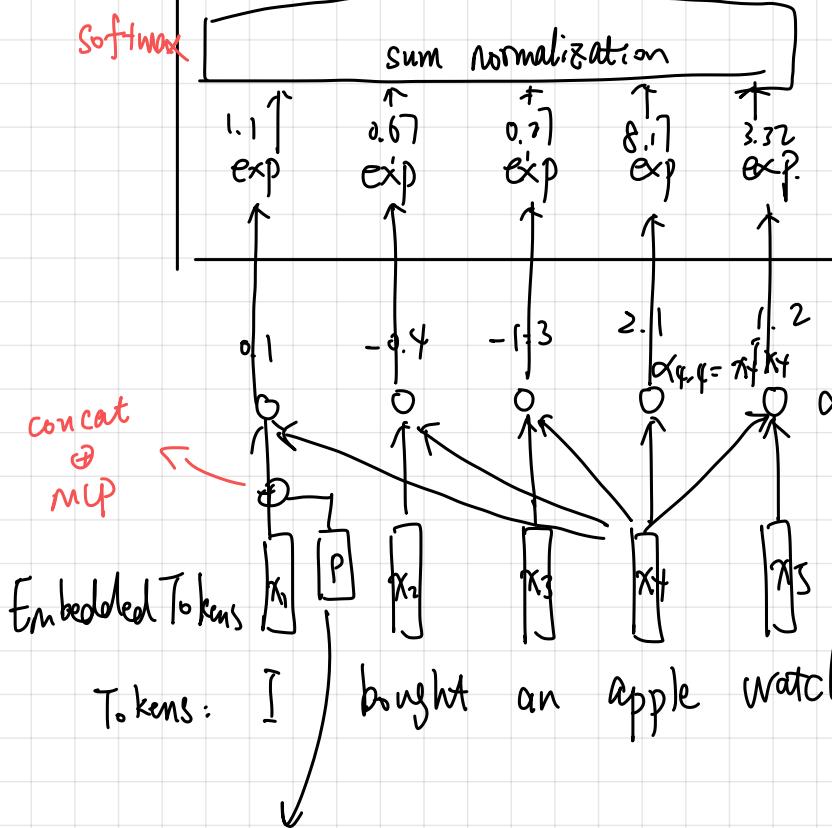
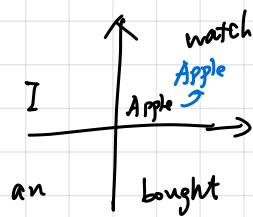
$$d \begin{bmatrix} 0.5 \\ 0.3 \\ 0.8 \end{bmatrix} = d \begin{bmatrix} \boxed{1} & \boxed{2} & \boxed{3} \end{bmatrix} \quad \begin{matrix} \text{How} \\ \# \text{token} \end{matrix} \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

How are you.
Embedded matrix.

But just embedded matrix is not enough,

Apple should rely on context to solve ambiguity

$$\Rightarrow \hat{x}_4' = \alpha'_{4,1}x_1 + \alpha'_{4,2}x_2 + \dots + \alpha'_{4,5}x_5$$



$\alpha'_{4,5} = x_4^T x_5$ dot product to calculate similarity

A concept relevance
is asymmetrical

delicious apple
 \times

The positional embedding derive intuition from binary encoding

Position 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15
 2^0 0 0 0 0 0 0 0 1 1 1 1 1 1 1 1 1 → low oscillation

d 2¹ 0 0 0 0 1 1 1 1 0 0 0 0 1 1 1 1

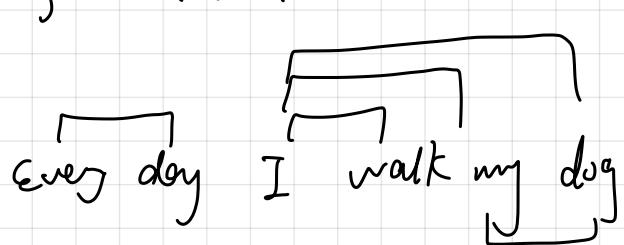
2² 0 0 1 1 0 0 1 1 0 0 1 1 0 0 1 1

2³ 0 1 0 1 0 1 0 1 0 1 0 1 0 1 0 1 → high oscillation

ViT

Relative positional encoding

I walk my dog every day \Leftrightarrow Every day I walk my dog

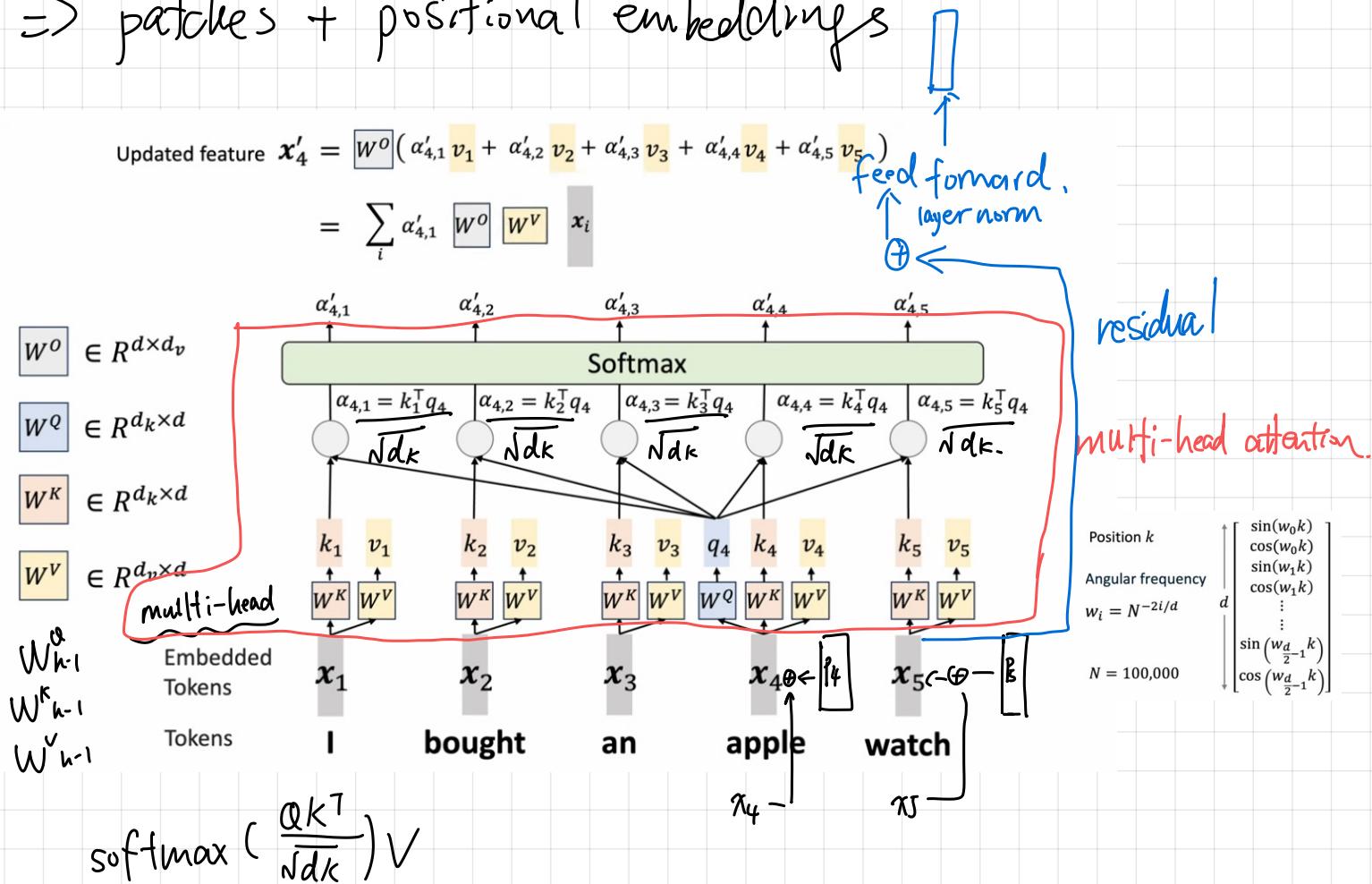


⇒ Rope

For pixels ?

Naive tokenization of a pixel for 224×224 image.
over 500k tokens ⇒ Too much

⇒ patches + positional embeddings



MAE.

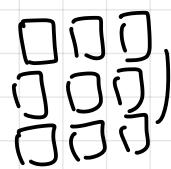
$$L = \sum (x_i - \hat{x}_i)^2$$

25% remaining patches
75% masked image patches

Encoder

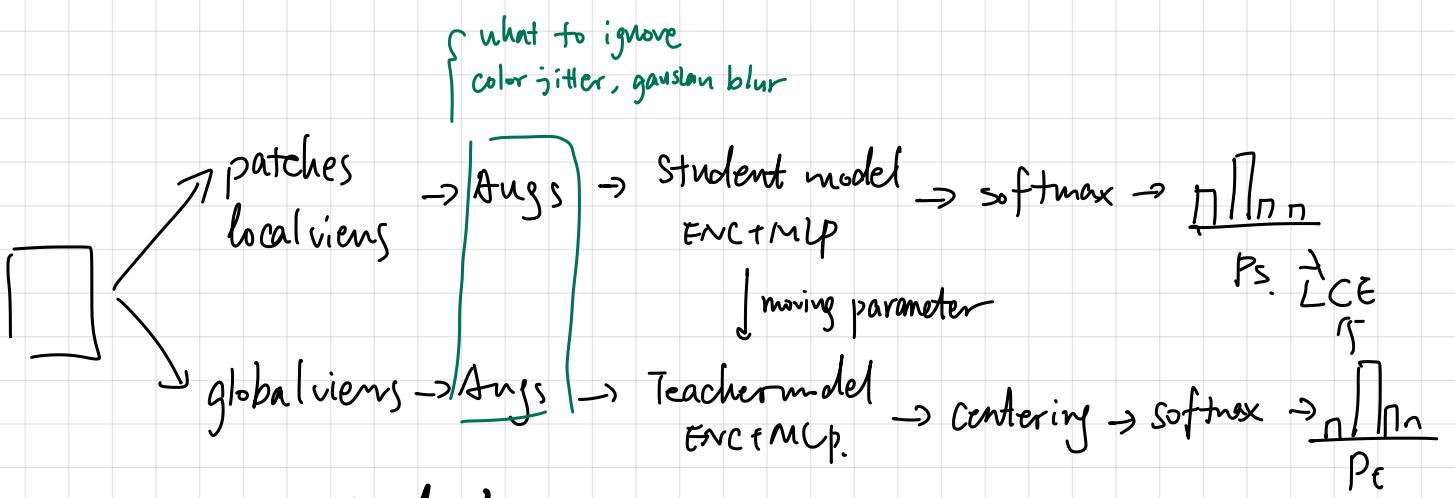
compact representation

Decoder

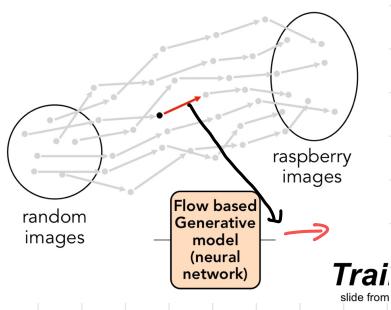


DINO Distillation with NO Labels

A student network learns from a teacher network



Diffusion Model



estimate the probability distribution of images $P(x)$
 ① sample from a gaussian distribution $\sim \mathcal{N}(0, I)$
 ② project to images

Algorithm 1 Training

```

1: repeat
2:    $\mathbf{x}_1 \sim$  clean image from training set
3:    $t \sim \text{Uniform}([0, 1])$ 
4:    $\mathbf{x}_0 \sim \mathcal{N}(0, I)$ 
5:    $\mathbf{x}_t = (1 - t)\mathbf{x}_0 + t\mathbf{x}_1$ 
6:   Take gradient descent step on
       $\nabla_{\theta} \|(\mathbf{x}_1 - \mathbf{x}_0) - u_{\theta}(\mathbf{x}_t, t)\|^2$ 
7: until happy
  
```

Algorithm 2 Sampling

```

1: input:  $T$  timesteps
2:  $\mathbf{x}_t = \mathbf{x}_0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
3: for  $t$  from 0 to  $T$ , step size  $\frac{1}{T}$  do
4:    $\mathbf{x}_t = \mathbf{x}_t + \frac{1}{T} u_{\theta}(\mathbf{x}_t, t)$ 
5: end for
6: return  $\mathbf{x}_t$ 
  
```

Algorithm 3 Class-Conditional Training

```

1: repeat
2:    $\mathbf{x}_1, c \sim$  clean image and label from training set
3:   Make  $c$  into a one-hot vector
4:   with probability  $p_{\text{uncond}}$  set  $c$  to zero-vector.
5:    $t \sim \text{Uniform}([0, 1])$ 
6:    $\mathbf{x}_0 \sim \mathcal{N}(0, I)$ 
7:    $\mathbf{x}_t = (1 - t)\mathbf{x}_0 + t\mathbf{x}_1$ 
8:   Take gradient descent step on
       $\nabla_{\theta} \|(\mathbf{x}_1 - \mathbf{x}_0) - u_{\theta}(\mathbf{x}_t, t, c)\|^2$ 
9: until happy
  
```

Algorithm 4 Class-Conditional Sampling

```

1: input: one-hot vector  $c$ , classifier guidance scale  $\gamma$ 
2:  $\mathbf{x}_t = \mathbf{x}_0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
3: for  $t$  from 0 to  $T$ , step size  $\frac{1}{T}$  do
4:    $u_{\text{uncond}} = u_{\theta}(\mathbf{x}_t, t, 0)$ 
5:    $u_{\text{cond}} = u_{\theta}(\mathbf{x}_t, t, c)$ 
6:    $u = u_{\text{uncond}} + \gamma(u_{\text{cond}} - u_{\text{uncond}})$ 
7:    $\mathbf{x}_t = \mathbf{x}_t + \frac{1}{T} u$ 
8: end for
9: return  $\mathbf{x}_t$ 
  
```

Objective: $\mathbb{E}_{z \sim q(z|x)} \log p(x_i|z) - DL(q(z|x_i)||p(z))$

$$\begin{aligned}
 D_{KL}(q_{\phi}(z|x_i)||p(z)) &= \mathbb{E}_{z \sim q_{\phi}(z|x_i)} \left[\log \frac{q_{\phi}(z|x_i)}{p(z)} \right] \\
 &= \mathbb{E}_{z \sim q_{\phi}(z|x_i)} [\log q_{\phi}(z|x_i)] - \mathbb{E}_{z \sim q_{\phi}(z|x_i)} [\log p(z)] \\
 &\stackrel{\text{Entropy}}{=} -\mathcal{H}(q_{\phi}(z|x_i)) \\
 -DL(q_{\phi}(z|x_i)||p(z)) &= \mathbb{E}_{z \sim q_{\phi}(z|x_i)} [\log p(z)] + \mathcal{H}(q_{\phi}(z|x_i))
 \end{aligned}$$

Re-written Objective: $\mathbb{E}_{z \sim q(z|x)} \log p(x_i|z) + \log p(z) + \mathcal{H}(q(z|x))$

Variational Autoencoders

- The Encoder maps x to z .
 - But need to be able to sample from z like $p(z)$
 - So we need the encoder to learn $p(z|x)$
 - But how to get $p(z|x)$...?
 - What's wrong?
 - So learn $q(z|x)$ and make it close to $p(z)$ $DL(q(z|x)||p(z))$
 - Through ELBO (will do later):
- $$\max \log p(x) \geq \mathbb{E}_{q(z|x)} \log p(x|z) - DL(q(z|x)||p(z))$$
- i.e. max $p(x|z)$ while regularize $q(z|x)$ to be close to $p(z)$