

Lecture 1: Basics

Capture Devices

Passive:

→ RGB

→ Stereo and Multi-View

+ work indoors & outdoors

- Need Computer

- Features needed

Active:

→ TOF (Time of flight)

→ Structured light

→ LiDAR

- slow

+ very precise

+ No features needed

- Fails outdoors

- Sensitive to battery

Stereo Matching: Sparse vs. Dense Matching

Wider dist. between cameras: Feature extraction vs. Nat.

↑ distortion → occlusion ↗

Smaller dist. between cameras: Search window: smaller = more noise
↑ area ↑

Same as TOF
but needs precise calibration between projector and sensor

- Correspondence Failures:
- 1) Non-Lambertian surfaces → specular reflection
(not diffuse)
 - 2) Occlusions / repetitions
 - 3) Textureless Surface

Camera Concepts

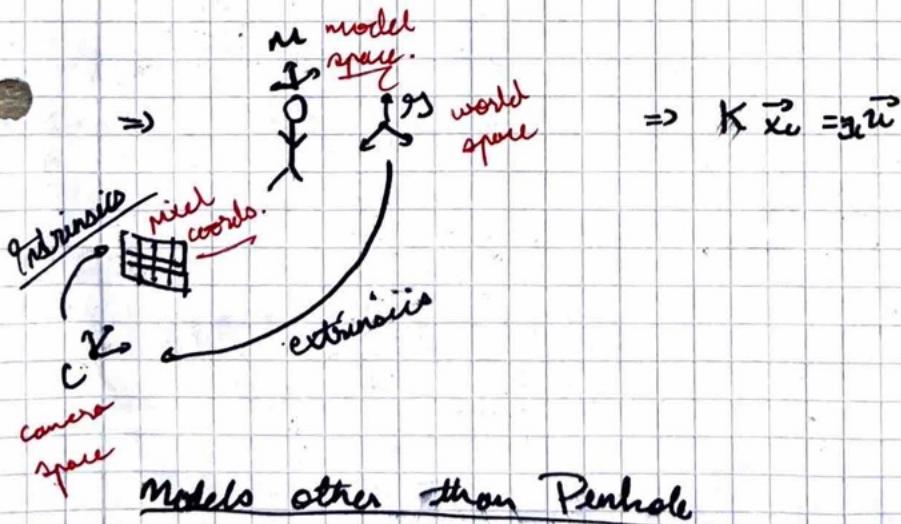
- Intrinsic: → Focal Length
→ Principal Point.
→ (Skew)
→ (Distortion)

$$K = \begin{pmatrix} f_x & r & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{pmatrix}$$

δ : shear coefficient
 c_x, c_y : principal point or optical center.

f_x, f_y : focal length.

- Extrinsic: R, t .



$$\Rightarrow K \vec{x}_w = \vec{y} \vec{u}$$

Models other than Pinhole

- Radial + Tangential Distortion:

Fish eye lens.

(large plane not parallel to lens)

use approximate values using camera calibration.

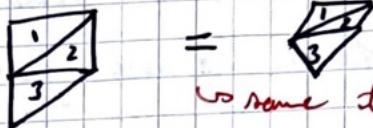
Lecture 2: Surface Representations

obviously or Polygonal Meshes:

- efficient
- secondary:
 - Error in shape representation: $O(l^2)$ ~ i.e. doubling the # of edges gives $\times \left(\frac{1}{2}\right)^2$ reduction in error.
 - Allows for recursive subdivisions (i.e. higher resolution)

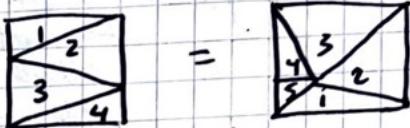
Topology (triangulation) vs. Geometry

Topologically equivalent



\Rightarrow same triangulation

Geometrically Different



\Rightarrow different triangulation

Geometrically equivalent

Topologically Different

\Rightarrow different topology but same rectangle.

Manifolds

→ Triangle is called manifold if:

mesh (i) Intersection of two triangles are (+ breve):

~~adjacent~~
~~shared~~
~~edge~~
~~but not vertex~~
~~Non-manifold~~

~~empty~~ → empty i.e. no area.
+ common vertex
and pl. + common edg.

here
intersection spans the area
of small triangle = not empty.
→ Non-manifold

(ii) Edges have:

- * One adjacent triangle → border edge
- * Two adjacent triangles → inner edge

(iii) Vertices:

- + Build a single open fan → border vertex
- + Build a single closed fan → inner vertex

edge
ab has three
adjacent triangle
 \Rightarrow Non-manifold

vertex connected to an open fan
two open fans
 \Rightarrow non-manifold

vertex is
connected to
two closed fans
 \Rightarrow Non-manifold

Mesh Data Structures

(i) Shared Vertex

+ OFF Obj STL etc. ~ PL4 allows for binary storage

Shared
Vertex
data
structure

+ list v → list f.
- - {vertices}
- - - space.

(ii) Half-Edge Data Structure

no only accepts manifold
representations (since stores
adjacent edges)

to store vertices + faces + edges (i.e. opposite edge
adjacent edge etc.)

Useful for geometric computations

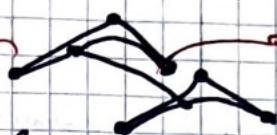
Surfaces

- (i) Explicit surfaces: $f: \mathbb{R}^2 \rightarrow \mathbb{R}$ for terrain in games
i.e. given x, y and 1 give height
- (ii) Parametric surfaces:

Constructed from Bezier Curves — Bezier patches

CAD, etc.
Used in 3D design
+ construction etc.

Bezier curve



Bezier patch:

$$f(u, v): \mathbb{R}^2 \rightarrow \mathbb{R}^3$$

- (iii) Constructive Solid Geometry:

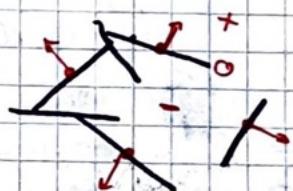
Boundaries created from boolean operations:



- (iv) Implicit Surfaces:
→ SDF
→ Height
→ Density

Defining Implicit Surfaces:

- local method
- an (i) Hoppe method: $f(\vec{x}) = (\vec{x} - \vec{p}) \cdot \vec{n}_p \approx \text{normal of point } \vec{p}$
for any \vec{x} → nearest \vec{p} (from point cloud e.g.)
- creates Voronoi diagram which defines closest points
- piecewise linear
- discontinuous along edges of voronoi



- (ii) RBF: Radial Basis Functions:

We assume smoothness $\Rightarrow f(\vec{x}) = \sum_{i=1}^n \alpha_i \varphi_i(\vec{x}) + b \cdot \vec{x} + d$

To solve: $f(\vec{p}_i) = 0$

$$f(\vec{p}_i + \epsilon \delta \cdot \vec{n}_p) = \epsilon \quad \left. \begin{array}{l} \text{our constraints} \\ \text{some C2 function} \end{array} \right\}$$

$$f(\vec{p}_i - \epsilon \vec{n}_p) = -\epsilon \quad \left. \begin{array}{l} \text{Unknowns: } \alpha_1, \alpha_2, \dots, \alpha_n, b, d \\ n+3 \end{array} \right\} = n+4$$

Linear \Rightarrow SVD to solve for α

- (iii) Poisson Reconstruction

$$\nabla \cdot \nabla f(\vec{x}) = \nabla \cdot \vec{N} \approx \text{normal vector field}$$

is basically divergence of all our normals

divergence of our $\nabla f(\vec{x})$

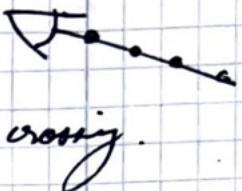
gradient field

for poison we don't use SDF

Rendering Implicit Surfaces

(i) Ray Casting

1) Ray Marching:
 - Apply fixed step length.
 - linearly Interpolate if zero crossing.

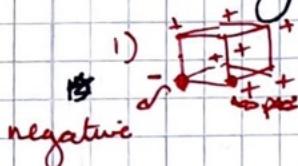


2) Sphere Tracing:
 - Dynamic Steps
 - Stop if distance below threshold.

(Also-to method
usually).



(ii) Marching Cubes: 1) Determine zero crossings for all grid cells.



2) If zero-crossing use Lookup table to triangulate

3) Linearly interpolate triangulation

negative
1)
2) Use which vertex negative as ID in Lookup Table.

3) ~~Find~~ Find Triangulation

4) Linearly interpolate triangulation

Lecture 3: 3D reconstruction Methods

Feature Detection

(i) FAST Detectors: for pixel $p(x,y)$ if ≥ 12 pixels are brighter

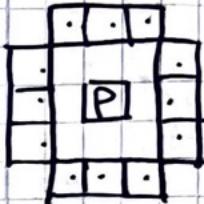
\Rightarrow feature point partial derivatives w.r.t. $x, y: I(u,v) - I(u+1,v)$ etc.

(ii) Harris Corner Detectors: if we have I_x, I_y :

we can do SSD for e.g. Harris but not robust.

$$H = \begin{pmatrix} \sum_w I_x^2 & \sum_w I_x I_y \\ \sum_w I_x I_y & \sum_w I_y^2 \end{pmatrix} \text{ is Hessian of some image patch } W. \text{ basically}$$

we find λ_1, λ_2 of H :
 $\lambda_1 \sim \lambda_2 + \text{small} \Rightarrow$ flat region
 $\lambda_1 > \lambda_2 \Rightarrow$ edge region
 $\lambda_1 \sim \lambda_2 + \text{large} \Rightarrow$ corner regions.



Trick to find λ_1, λ_2 w/o SVD: $R = \det(M) - \alpha (\text{trace}(M))^2$
 $= \lambda_1 \lambda_2 - \alpha (\lambda_1 + \lambda_2)^2 \Rightarrow$ properties of det + trace

if $\lambda_1 \sim \lambda_2 \sim 0 \Rightarrow R \sim 0 \Rightarrow$ i.e. flat region
if $\lambda_1 > \lambda_2 \Rightarrow \lambda_1 \lambda_2 \sim 0 \Rightarrow R = -\alpha \lambda_1^2 \Rightarrow$ edge
if $\lambda_1 \sim \lambda_2 \gg \lambda_1 \lambda_2 \Rightarrow -\alpha \lambda_1^2 - \alpha \lambda_2^2 \Rightarrow$ corner.
is usually $0.04-0.06$

or if: $|R|$ small \Rightarrow flat region
 $R < 0 \Rightarrow$ edge regions
 $R > 0 \Rightarrow$ corner regions

{ cool trick actually :)

(iii) SIFT Detectors: finds features that are scale invariant + filter (usually Haar LS)
(Cand rotation invariant) i.e. applies multiple scalings + blur and finds features not affected + finds dominant rotation direction to compare features.

Correspondence Matching

(i) Simple SSD per feature: shitty + not robust

(ii) SIFT Descriptors: after SIFT filters optimal features: takes 16-sub blocks neighbourhood around feature point, rotates neighbourhood in preferential direction, calculates 8-bin orientation histogram per sub-block $\Rightarrow 8 \cdot 16 = 128$ dimensional feature vector per feature / guess.

(iii) ConvNet: Basically siamese (shared weights) network inputting some image/volume patch from two views and outputs 1/0 depending on match or not.

(iv) SURF, ORB etc.

Main 3D recon methods:

(i) Visual Hull Carving: If we have:

BRUNNEN

- \rightarrow 2 cameras
- \rightarrow we know each camera's intrinsics
- \rightarrow we know each camera's extrinsics
- \rightarrow we segment subject in each camera

Carve over voxel grid based on segmentation

←

(ii) Structure from Motion (SfM): Given:

- Set of images
- Features + correspondences for each image

Bundle Adjustment

Re-projecting energy.

$$\text{extinguishes over time} \rightarrow \text{Res-proj}(T, X) =$$

$$x_j^* T_i +$$

Images i

$$\sum_{i=1}^m \sum_{j=1}^{n_i} \|x_{ij} - \pi_i(T_i x_j)\|_2^2$$

Bundle Adjustment

Goal: Find 3D locations X_j
+
Find Camera poses T_i .

feature points j of feature point j to camera frame i : $T_i x_j$
in image i then projected to camera i .

- * Unknowns: → $6 \cdot (m-1)$ → Extrinsics of all cameras except first (axis identity)
→ $3 \cdot n$ → Feature points in 3D space
→ ~~Camera optional~~ → Camera Intrinsic per camera.

- * Constraints: → $2 \cdot m \cdot j$ from all feature points in all images in 2D space.

$$\Rightarrow 2 \cdot m \cdot j \geq 6(m-1) + 3j + k.$$

(iii) Multiple View Stereo (MVS): Given:

- We can do frame pair matching + vote

→ Poses of images T_i
→ 3D locations of features X_j

Goal: → Dense Point Cloud

- Global optimization is NP-Hard

- Plug point clouds to surface fitting e.g. person surface fitting

Usually Pipeline: 2D images → SfM → MVS → Reconstruction → 3D mesh.

(iv) SLAM: ~ smaller compute than GM + MVS.

Frame to frame OR

Frame to Model

= Basically real-time SfM + MVS

≠ Different than Bundle Adjustment since no simultaneous offline global optimization of all poses + 3D locations.

* Types: { 1) ORB SLAM } can also either use feature points for alignment OR directly input
RGB { 2) LSD SLAM. }

RGB-D { 3) KinectFusion }

SLAM usually has Loop-Closure Problem

RGB-D Bundling

$$\text{Ebundle}(T) = \sum_{i,j} \sum_n \|T_i p_{in} - T_j p_{jn}\|_2^2$$

$$\text{Edepth}(T) = \sum_{i,j} \sum_n \|(\rho_n - T_i^{-1} T_j^{-1} \pi_j(D_j(\pi_j(T_j T_i p_{in})))) \cdot n_n\|_2$$

$$\text{Ecolor}(T) = \sum_{i,j} \sum_n \|\nabla I(\pi_i(p_{in})) - \nabla I(\pi_j(T_j^{-1} T_i p_{in}))\|_2$$

Color gradient of original pixel in frame i

Color gradient of projected pixel in frame j

Lecture 4: Optimization Methods:

Linear Least Squares

→ Means the model is linear in parameters

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 \dots$$

↳ here X_2 can be x^2 , $\log x$, $\frac{1}{x}$ doesn't matter
still linear since linearly multiplying by β_2 .

→ Simplified to: $Ax = b$

↳ overdetermined usually

→ Usually solved as normal equation: $A^T A x = A^T b$

Solvers: vs of linear least squares

→ Iterative Solvers: conjugate

- 1) Preconditioned Gradient Descent. → most common iterative solver.
- 2) Gauss-Seidel Iteration
- 3) Jacobi Iteration

→ Direct Solvers:

- related to each other I think
- 1) QR, LU - Decomposition
 - 2) SVD
 - 3) Cholesky

Non-Linear Least Squares

* No longer linear in params, e.g. $Y = \beta_0 X^{\beta_1}$ ↳ β_1 not linearly multiplied

Solvers: we're trying to solve $\underset{x}{\operatorname{argmin}} f(x) = \|F(x)\|_2^2$

↳ First Order Solvers: Gradient Descent etc.

$$\Rightarrow x_{k+1} = x_k - t \cdot \nabla f(x_k) \text{ ↳ momentum / line search standard BS.}$$

→ Second Order Solvers:

(i) Newton's Method: $x_{k+1} = x_k - H_f(x_k)^{-1} \nabla f(x_k)$ ↳ comes from Taylor

$$\text{taylor: } f(x+t) \approx f(x) + \nabla f \cdot t + \frac{1}{2} H_f \cdot t^2$$

(ii) Gauss-Newton: Approximation of Newton:

$$\underbrace{H_f \approx 2 J_F^T J_F}_{\substack{n \times n \\ \text{residuals where } F \in \mathbb{R}^n}} \Rightarrow x_{k+1} = x_k - [2 J_F^T J_F]^{-1} \nabla f(x_k)$$

$$\Rightarrow f(x+t) = \nabla f + H_f \cdot t = 0 \text{ ↳ } t = H_f^{-1} \nabla f. \text{ so: } x_{k+1} = x_k - t.$$

$$\underbrace{2 J_F^T J_F}_{\substack{\text{linear system } \text{tow. } \text{P.P.} \\ \text{A}}} \underbrace{[x_k - x_{k+1}]}_x = \underbrace{\nabla f(x_k)}_b \quad (\text{PCG})$$

→ solve with iterative methods like preconditioned gradient descent linear solvers.

(iii) Levenberg: $x_{k+1} = x_k - [2 J_F^T J_F + \lambda I] \nabla f(x_k)$ ↳ conjugate

$$\lambda \rightarrow \infty \text{ N}$$

$\lambda \uparrow \rightarrow$ Gradient Descent.

↳ Tikhonov regularization

i.e. damping factor for each term in J s.t. $\lambda \uparrow$ if $f(x_k) < f(x_{k+1})$

(iv) Levenberg - Marquardt (LM): faster convergence this way than Levenberg

$$x_{k+1} = x_k - [2\mathbf{J}_F^T \mathbf{J}_F + \lambda \cdot \text{diag}(\mathbf{J}_F^T \mathbf{J}_F)]^{-1} \nabla f(x_k)$$

more sophisticated Levenberg since each λ scaled by curvature of (diagonal is independent parts of curve)

(v) BFGS / LBFGS:

$$x_{k+1} = x_k - B_k \nabla f(x_k) \quad \text{Basically approximation of Hessian.}$$

$$\text{where } B_{k+1} = B_k + \alpha u u^T + \beta v v^T$$

In Practice we use Limited Memory BFGS (LBFGS)

Outlier Handling:

? How do we convexify our problems?
course to fit over residuals

(i) ~~RANSAC~~: basically trial & error

(ii) Lifting Schemes:

$$\text{gives: } f(x) = \sum r_i(x)^2 \Rightarrow f_{\text{robust}}(x) = \sum w_i^2 r_i(x)^2 + \lambda \log \sum (1-w_i^2)^{-1}$$

(iii) Robust Norms: usually hard to optimize | that penalizes $w \rightarrow 0$
(intials are $w=1$)

\Rightarrow Iteratively Reweighted Least Squares

$$f_{\text{robust}}(x) = \sum w_i r_i(x)^2 \text{ and } w_i = |r_i(x)|^{p-2} \text{ for some } p\text{-norm.}$$

$$\text{assuming we wanted } \underset{x}{\operatorname{argmin}} f(x) = \underset{x}{\operatorname{argmin}} \|F(x)\|_p^p$$

Derivative Types:

(i) Numeric derivatives: $\frac{df(x)}{dx} \approx \frac{f(x+h) - f(x)}{h}$

\rightarrow Easy to implement \Rightarrow good for debugging

\rightarrow Slow and numerically unstable

(ii) Automatic Differentiation: Use dual numbers: $e \neq 0$

$$\text{Example: } f(x) = x^2 \quad \frac{df(x)}{dx} \Big|_{x=10} ? \Rightarrow f(10+e) = 100 + 20e + e^2 \quad \begin{matrix} e^2 = 0 \\ \cancel{e^2} = 0 \end{matrix} \quad \Rightarrow \frac{df}{dx} \Big|_{x=10} = \frac{20}{\cancel{e}}$$

(iii) Symbolic Differentiation:

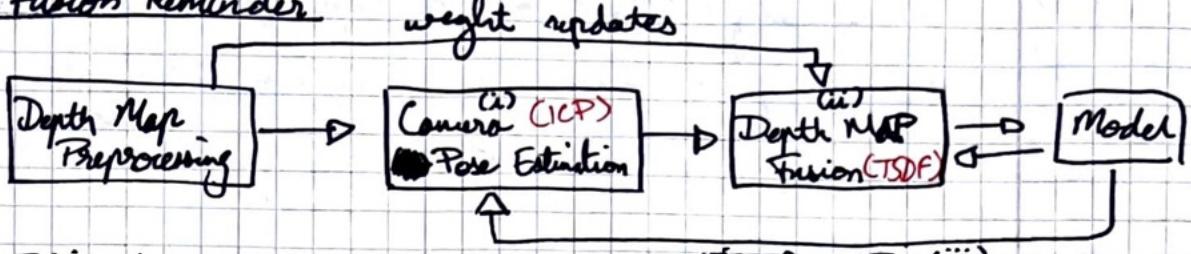
Basically no numerical solutions (Maple :))

Other Solvers: 1) Inequality Constraints: Lagrange multipliers etc.

2) Gradient-Free: Monte Carlo etc.

Lecture 5: Rigid Surface Tracking + Recon (i.e. KinectFusion revision)

KinectFusion Reminder



(i) Pose Estimation via ICP:

Basic Steps of ICP:

! If coarse to fine registration for concave hull problem -
iterate

- 1) Selecting source points \rightarrow
- 2) Find Correspondences
- 3) Weigh Correspondences
- 4) Reject Outliers
- 5) Assign Errors metric
- 6) Minimize errors for pose estimate

Either:
 → Use all points
 → Uniform subsampling
 → Random sampling
 → Stable sampling

stable:

uniform:

random:

Finding Correspondences in ICP:

only works if we have depth values

a) Nearest neighbour \rightarrow FLANN etc.

b) Projective correspondences \rightarrow did in the project i.e.

* We can improve matching via only checking compatible points:

- e.g. at border.
- a) Color Compatibility
 - b) Normal Compatibility
 - c) Pruning

- d) Reject angle and distance outliers

~~ps $\rightarrow \pi(T(T^{-1}(T(p)))$)~~
 my correspondence in other frame.

Errors in ICP:

had closed form linear solution??

idea here

is energy

lowest if

$p_i \in q_i$ on

a) Point to Point: $\sum_i \|(\mathbf{T}p_i - q_i)\|_2^2 = E(T)$

b) Point to Plane: $\sum_i \|(\mathbf{T}p_i - q_i) \cdot n_i\|_2^2 = E(T)$

usually solved w/ LM
 if not linearized

\hookrightarrow can be linearized in 2D
 assuming small rotations \Rightarrow cost \propto $\|\mathbf{Cp} - \mathbf{q}\|$
 \Rightarrow Cholesky / SVD

* Symmetry ICP: $C_p - q \propto (n_p \pm n_q)$

(ii) SDF / TSDF: (Implicit Representations):

- Basically representation of distance to zero-crossing band.

- TSDF: just applies truncation if > 1 from zero-crossing

$/ < -1 \Rightarrow 1 / -1$ everywhere

\hookrightarrow Surfel representation: point + normal + radius

Volumetric Fusion: also exists current tsdf based on our current depth value

weights of running average concept: $D_{i+1}(x) = W_i(x)D_i(x) + W_{i+1}(x)d_{i+1}(x)$

can either be 1

everywhere OR closer

to camera weight, away weight (unweighted)

tsdf value at voxel i currently (i.e. distance to no surface) $W_{i+1} = W_i(x) + w_{i+1}(x)$

(iii) Surface Rendering

Important to
rely on SDFs:

=0 → surface
<0 → unknown
>0 → free space.

(could have another
surface ~~behind~~ behind)

a) Raycasting: Throw multiple rays from camera.

→ Surface found from sign change + interpolate to find intersection

→ Normal found from VTSDR.

b) Marching Cubes: Basically use Marching Cubes table
based on zero-crossing.

Follow-Ups to KinectFusion:

1) Extended Kinect Fusion:
+ easy to implement
+ unlimited spatial volume.
- cannot re-integrate previously regions.

2) Hierarchical Fusion:
(basically uses sparse voxel octrees)
+ fast lookups
+ efficient storage (high local res)
- costly updates
(insertions/removals)

3) Voxel Hashing:
+ extremely fast updates (insertions/removals) O(1)
+ efficient storage
- lookups slower => raycasting not as fast
+ cheap streaming

Loop Closure Follow-Ups:

1) Elastic Fusion:
→ Surfel based
→ Detects loop closure by localization
→ Non-rigidly warp to close loop.

2) Bundle Fusion:
→ Added sparse energy term for pose estimation
→ Includes surface de-integration if pose had wrong estimate

Lecture 6: Deformations of Non-rigid Surface Tracing

Mesh Deformations → standard

$$E = E_{\text{fit}}(V) + E_{\text{reg}}(V, X)$$

we pick non-rigid regularizers which affects behaviour of our mesh.

$$E_{\text{fit}}(V) = \sum_{i \in C} \|c_i - v_i\|_2^2$$

minimize distance between vertex and target

1-ring fan: Regularizer Types → Our goal for regularizers: → Smooth Global shape → Global changes preserve local detail → Should be zero in undeformed state.

i) Laplacian Surface Editing:

$$E_{\text{reg}}(V, X) = \sum_{j \in N} \|(v_i' - v_i) - (v_j' - v_j)\|_2^2$$

1-ring fan neighborhood (1-ring) → i.e. deformation of v_j' should be very similar to 1-ring neighbours v_i'

ii) As Rigid as Possible (ARAP):

$$E_{\text{reg}}(R, V) = \sum_i \sum_{j \in N} \|(v_i' - v_j) - R_i(v_i - v_j)\|_2^2$$

done per fan over all vertices means that after some R edge

* we could add weight to give bigger priority to closer neighbours: $w_{ij} = \frac{1}{2}(\cot \alpha_{ij} + \cot \beta_{ij})$

$$\frac{v_i' - v_j'}{\|v_i' - v_j'\|} = \frac{v_i - v_j}{\|v_i - v_j\|}$$

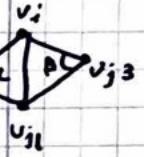
v_i → v_i' after some rotation.

* Solving: Unknowns: $3n + 3n$ per vertex

$$\text{Constraints: } 3n \text{ valence} + 3c$$

average a) from E_{fit}
per position
for neighbours

so $\cot \alpha_{ij} + \cot \beta_{ij}$
⇒ if v_{ij} closer then v_{ij} means larger



→ We can do Levenberg Marquardt (non-linear)

OR

→ Flip-flop optimization: 1) Assume v_i', v_j' constant and solve for R

2) Assume R constant and solve for v_i', v_j'

⇒ linear problem and can be solved with SVD / Procrustes.

iii) Embedded Deformation: reduced version of ARAP

~~$$E_{\text{reg}}(M, V) = \sum_{j \in N} \|(v_i' - v_j) - M_i(v_i - v_j)\|_2^2$$~~

$$\text{Error} = \sum_i \text{Rot}(M_i), \quad \text{Rot}(M) = (c_1 \cdot c_2)^2 + (c_1 \cdot c_3)^2 + \dots + (c_1 \cdot c_n)^2$$

* Solving: Unknowns: $12n$

(i.e. we want M to resemble R but allow for more DoF in our deformation (shear etc.))

⇒ quartic problem to optimize

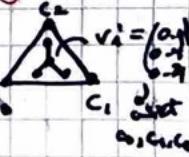
Deformation Principles instead of deforming mesh = computationally expensive deform proxy

we can apply: i) Cage: defines combination of vertices w.r.t. cage.

(we use barycentric coordinates for that.)

ARAP / ED
Implementation here.

ISRUNNEN ⇒ we deform cages and by proxy our vertices using barycentric coordinates.



ii) 3D Grids:

e.g. blend skinning or
iii) Skeletons: hierarchical:
Body
left shoulder
head.
hand

Called
linear
blend skinning

Each vertex where there is a linear combination of bone
matrices: $v_i = \sum w_{ik} M_k v_i'$
weight for each bone has a matrix
each bone

iv) Deformation Graphs: assigns a region/neighborhood a control
node. is similar to cage but instead we have deformation
zone.

Non-Rigid Tracking is full correspondences than ICP.

needs pre-defined
mesh.

* Non-Rigid ICP: more degrees of freedom. e.g. ARAP/ED.

$E = \lambda \text{point Epoint} + \lambda \text{plane Eplane} + \lambda \text{color Eccolors} + \lambda \text{key Ekey}$.

from regular point to plane
ICP

i.e. lower
energy if
colors match

$$\text{Epoint} = \sum_i \|v_i - T v_i'\|^2$$

$$\text{Eplane} = \sum_i [(v_i - T v_i') \cdot n v_i]^2 \text{ forces } v_i' \text{ to stay on the plane}$$

$$\text{Eccolors} = \sum_i [I(\pi(v_i)) - I'(\pi(T v_i))]^2 \text{ of } v_i$$

$$\min_{T, v_i'} \|v_i - T v_i'\| \Rightarrow E=0$$

Ekey: ARAP. lower energy if colors match.

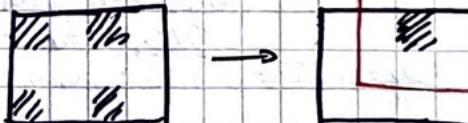
* Important hyperparameter-thresholds for non-rigid ICP:

- 1) distance
- 2) angle
- 3) Occlusion
- 4) Viewing angle

Improves with more cameras. for speed.

* Correspondence Association:

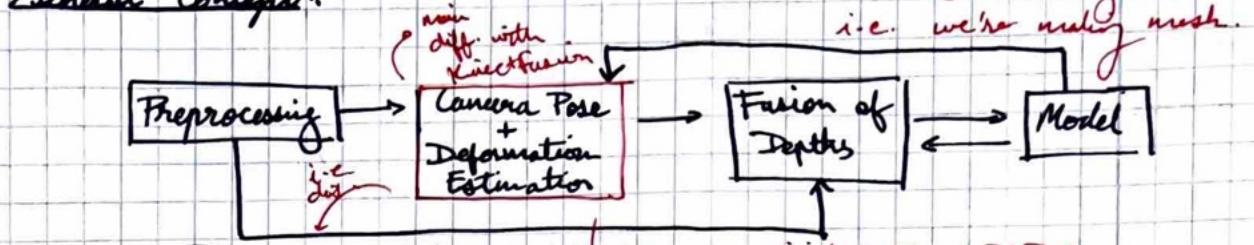
c.g.:
1) Projective correspondences every k-th element
2) Kernel around k-th sample for minimum reduction



to find best
match within kernel.

Lecture 7: Non-Rigid Tracking and Reconstruction

General Concepts:



* For Camera Pose + Deformation Estimation i.e. non-rigid ICP + ARAP etc.

- iterate
 - 1) Deform SDF
 - 2) Find isosurface or usually with marching cubes since raycasting hard for
 - 3) Find Correspondences
 - 4) Optimize deformations non-rigid ICP with some non-rigid regularizers
 - 5) Integrate RGB-D data to model

fast motion
ICP passes

Main Methods discussed:

th reconstruction: (i) DynamiC Fusion: ↗ Uses Deformation Graphs (ED) ↗ than this.

course is

surface moving ↗ Main OK skin deforming? ↗ We can increase resolution based on grid ↗ This means higher stability compared to dynamic fusion ↗ uses regular deformation grids + ARAP ↗ Non-rigid ICP using sparse SIFT matches ↗ Volumetric Fusion ↗ More stable + Better texture.

(ii) Volume Deform:

- Deformation Graphs (ED)
- Non-rigid ICP (dense depth + dense colors) instead of SIFT.
- illumination correction

(iii) Fusion 4D:

- way more complex.
- 8 Depth Maps (16 IR cameras)
- Key volumes: multiple poses → i.e. allows for topological changes guy running coat etc.

(iv) Motion2Fusion:

- Basically faster Fusion 4D: 200FPS

→ ED Graphs

→ Key Volumes

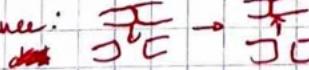
→ Multiple Depth Cameras.

→ i.e. solving topology problem.

* Motion2Fusion had 2 novelties:

- 1) 2 way backward and forward solves it by matching ↗ nonrigid alignment

reference to date then (date to reference):



(1) close to open; open to close:

reference to date to data.

- 2) Detail layer: they increased resolution of grid at surface (viewing rays)

(v) Lookin Good: ↗ increase / improve reconstructions using deep architecture

(vi) DeepDeform: ↗ find better correspondences than SIFT for best tracking

The two RGBs ↗ are the source ↗ query points and target where correspondences need to be found.

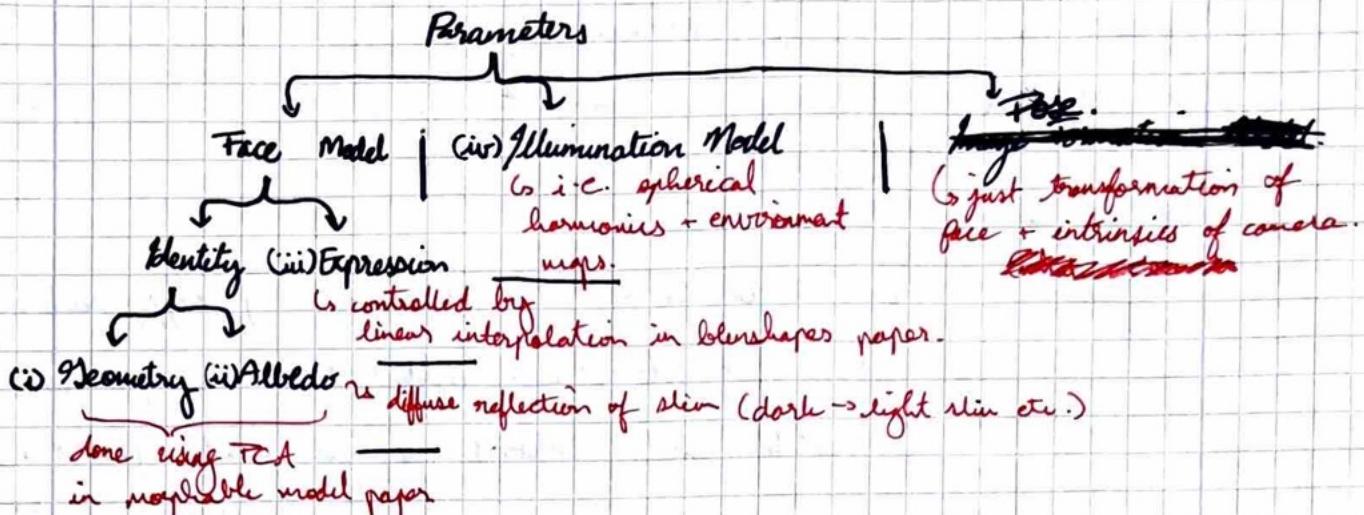
→ use siamese network that shares weights for two RGBs to give you correspondences.

$\star \Rightarrow$ For non-rigid reconstruction: no accumulates drift if tracking is wrong.

Tracking with a)

pre-defined mesh mitigates drift w.r.t. base mesh.

Lecture 8: Parametric Face Models



(i) Identity: Geometry

* Here by geometry we mean ~~the~~ structure of face (can include expressions)

* Steps ~~1 to 4~~: According to face warehouse + morphable model

1) Scan a bunch of faces

2) Fit topologically consistent template non-rigidly to faces

↳ Since all registration from same model we have 1 to 1 ~~correspondences~~ correspondences

⇒ 3) Create one massive vector of all positions of vertices per face

4) Find PCA-basis for all the vectors (i.e. all faces)

⇒ $M_{\text{geo}}(\alpha) = \text{mid} + E_{\text{id}} \alpha$ ↳ linear parameter.
average shape ↳ PCA basis that we found.
PCA side note: ↳ this also means we get average face like this.

~~PCA~~

eigenvectors of the

* The principle components are the n largest eigenvalues of covariance matrix

↳ also means direction of largest variance in our data or ~~line~~ line with minimal dist. to data.

Covariance:

$$\text{cov}(X, X) = \frac{1}{N-1} \sum_i \bar{X}_i \bar{X}_i^T$$

where $\bar{X}_i = X_i - \bar{X}$ is average X_i

(ii) Identity: Albedo

* Here we mean light/dark skin

* Steps exactly the same as Geometry but done independently of Geometry.

~~PCA~~

average albedo ↳ PCA basis we found

$$M_{\text{alb}}(\beta) = \alpha_{\text{alb}} + E_{\text{alb}} \beta \quad \text{↳ linear parameter.}$$

(iii) Expression:

- * We can parametrize expressions by:

done usually to 1) PCA no semantic meaning though:
2) BlendShapes ^{↳ has semantic meaning}

- * BlendShapes: Steps:

- 1) Artist makes n number of target expressions (e.g. happy sad etc.)
- 2) We give a weight (blendshape weights) to each expression
then sum/interpolate all expressions to get final expression
(e.g. 0.5 happy + 0.5 sad.)

(iv) Illumination Model:

Steps

- Environment Maps
- Spherical Harmonics → parametrization of environment maps
 - 1) Environment maps place the object a cube/sphere map of some panoramic scene
 - 2) Project lighting onto object (which is in center of sphere/cube)
 - Assumptions: → Distant light
→ No scattering
- Spherical harmonics parametrize/reduce dimensionality of environment maps. (orthogonal basis)
 - ⇒ we can linearly interpolate basis spheres for different lighting
 - Assumptions: → Lambertian Surface:
→ Distant Smooth lighting
 - i.e. diffuse lighting depends on n^2

Lecture 9: Face Tracking and Reconstruction

* Goal: find parameters \mathbf{p}^* of parametric face model to fit RGB 2D image.

* We make energy term: \rightarrow for energy minimization

$$E(\mathbf{P}) = E_{\text{dense}}(\mathbf{P}) + E_{\text{sparse}}(\mathbf{P}) + E_{\text{reg}}(\mathbf{P})$$

Energy of
color and geometry of
mesh wrt. \mathbf{P}

Correspondence
terms wrt. \mathbf{P}

regularization of
 \mathbf{P}

$$E_{\text{dense}}(\mathbf{P}) = E_{\text{geom}}(\mathbf{P}) + E_{\text{col}}(\mathbf{P})$$

point to point \hookrightarrow color.
point to plane

point to plane

$$E_{\text{geom}}(\mathbf{P}) = \sum_i \| \mathbf{M}_i(\mathbf{P}) - D(\pi(\mathbf{M}_i(\mathbf{P}))) \|_2^2 + [(\mathbf{M}_i(\mathbf{P}) - D(\pi(\mathbf{M}_i(\mathbf{P})))) \cdot \mathbf{n}_i]^2$$

squared over all vertices \hookrightarrow mesh vertex

correspondent depth
value at projected
mesh vertex in
camera frame

dot product
of mesh mesh normal

$$E_{\text{col}}(\mathbf{P}) = \sum_i \| \mathbf{M}_i(\mathbf{P})_c - I(\pi(\mathbf{M}_i(\mathbf{P}))) \|_2^2$$

colors of
mesh at vertex i
(in mesh)

color of
pixel at projected
mesh vertex in camera
frame.

$$\left| \begin{array}{l} E_{\text{reg}}(\mathbf{P}) = \sum_k \left(\frac{\partial \mathbf{a}_k}{\partial P_{ik}, k} \right)^2 + \sum_l \left(\frac{\partial \mathbf{b}_l}{\partial P_{il}, l} \right)^2 \\ + \sum_m \left(\frac{\partial \mathbf{s}_m}{\partial P_{im}, m} \right)^2 \end{array} \right.$$

$\mathbf{a}, \mathbf{b}, \mathbf{s}$ are parameter
vectors of id, alr, exp
of centered at zero ofc.
The summands done for
each dimension

$$E_{\text{sparse}}(\mathbf{P}) = \sum_j \| \pi(\mathbf{M}_j(\mathbf{P}))_m - C_j \|_2^2$$

2D coordinates of projected
mesh in camera frame

is coordinates of correspondences
in 2D image

$\Rightarrow \mathbf{P}^* = \underset{\mathbf{P}}{\operatorname{arg\min}} E(\mathbf{P})$, solved with non-linear least
squares.

\hookrightarrow Gauss Newton or LM solver.

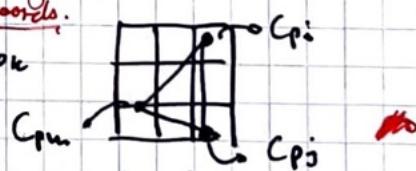
$$\Rightarrow \mathbf{P}_{n+1} = \mathbf{P}_n - (\mathbf{J}_F^T \mathbf{J}_F)^{-1} \mathbf{J}_F^T \cdot \mathbf{F} \quad \sim \text{Usually done Rough} \rightarrow \cancel{\text{Dense}} \text{ Dense pyramid}$$

Partial Derivatives of 2D image wrt. \mathbf{P} :

1) For each color pixel:

$$C_{\text{pixel}} = \alpha C_{pi} + \beta C_{pj} + \gamma C_{pk}$$

linear interpolation
of vertex color.



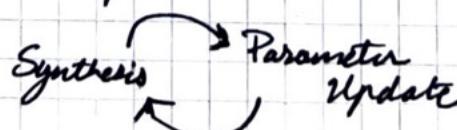
\hookrightarrow needed for Ecolor + Etex

\hookrightarrow probably the
same but not
sure :)

2) record triangle v , barycentric coords α, β, γ .

$$3) \Rightarrow \frac{\partial C_{\text{pixel}}(\mathbf{P})}{\partial P_j} = \alpha \frac{\partial C_{pi}(\mathbf{P})}{\partial P_j} + \beta \frac{\partial C_{pj}(\mathbf{P})}{\partial P_j} + \gamma \frac{\partial C_{pk}(\mathbf{P})}{\partial P_j}$$

\hookrightarrow This parameter update done:



* Can be done without depth via multiple views

(\hookrightarrow energy jointly optimized
for 6 images)

* Energy optimization dependent on:

→ Input Data (RGB vs. RGB-D)

→ Face Model + Illumination (coarse vs. dense)

→ Run-time requirements, hardware

→ Application

* Two main papers using this:

→ Real-time expression transfer for facial reenactment
 \hookrightarrow for two live videos

→ Face2Face

\hookrightarrow expression to RGB video

Lecture 10: Body & Hand Tracking

- Parameterization for consistent topology of hand + body
- Parameters of body: $M(\vec{\omega}, \vec{p}, \vec{u} | \phi)$
 - pose of body
 - shape
 - texture
 - hyperparameters.

CAESAR Dataset

→ 125 m 125 F

→ 74 markers.

→ Various ethnicities, age, weight.

- Same idea as before: average mesh + Non-rigid ICP
2DPCA of all mesh vectors.

SCAPE:

→ joint shape + pose learning.

→ blend shapes used for poses.

Skeleton Side-note:

$$\begin{array}{c} \vec{\omega}_1 \quad \vec{\omega}_2 \\ \text{Let } j_1 \quad j_2 \\ \vec{t}_s \end{array} \quad P_B : \text{each joint has a transformation matrix related to rotation by Rodrigues formula}$$

$$\Rightarrow \vec{P}_s = gS(\vec{\omega}_1, j_1) gS(\vec{\omega}_2, j_2) P_B$$

→ poses are given by $\theta = (\vec{\omega}_1, \dots, \vec{\omega}_n)$ where $\vec{f} = (j_1, \dots, j_n)$

Linear Blend Skinning: For any vertex \vec{t}_i on mesh:

vertices are linear combinations of all joints $\vec{t}_i' = \sum_{k=1}^n w_{k,i} g_{k,i}'(\theta, \vec{f}) \vec{t}_i$

weight given to transformation of each joint

Corrective Blend Shapes: each joint

Basically displace each vertex by $P = \begin{bmatrix} \Delta x_1, \Delta y_1, \Delta z_1 \\ \vdots \\ \Delta x_n, \Delta y_n, \Delta z_n \end{bmatrix}$
to simulate muscle compression.)

Non-rigid registration: $E = w_d E_d + w_s E_s + w_m E_m$

$$E_d = \sum_i w_i \text{dist}^2(T_i, T_i')$$

Data error.

$$E_s = \sum_i \|T_i - T_j\|_F^2$$

Transformation Smoothness

$$E_m = \sum_i \|T_i v_{ki} - m_i\|^2$$

marker Errors.

Hand Models:

Taylor: Uses subdivision model \rightarrow recursive coarser mesh to get:

→ C_2 continuity

→ Fast (1M) convergence

Issues with Model Generation + Fitting

- Expressive Scanning setup
- Neutral Pose similar but not identical
- Ambiguities in fitting \Rightarrow shape or expression?
- Drift in non-rigid registration