

## Lecture 1: Basics

### Capture Devices

#### Failures:

→ 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 lighting

### Stereo Matching: Sparse vs. Dense Matching

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

1 distortion  $\rightarrow$  occlusion  $\leftarrow$

Smaller dist. between cameras: Search window: smaller = more noise  
larger = less details

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

## Camera Concepts

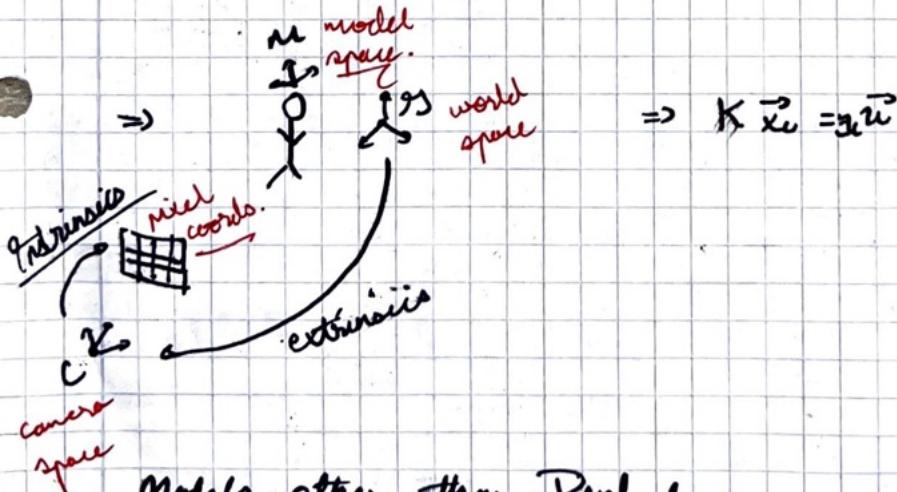
- Intrinsic: → Focal Length  
→ Principal Point.  
→ (Skew)  
→ (Distortion)

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

$\delta$ : shear coefficient  
 $c_x, c_y$ : principal  
points or  
optical center.

$f_x, f_y$ : focal  
length.

- Extrinsic:  $R, t$ .



### Models other than Pinhole

- Radial + Tangential Distortions: Use approximate values using camera calibrations.

Fish eye lens.

( $\rightarrow$  image plane not parallel to lens)

## Lecture 2: Surface Representations

obviously on Polygonal Meshes:

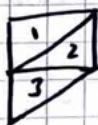
efficient

memory

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

### Topology (triangulation) vs. Geometry

Topologically equivalent



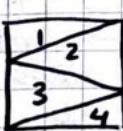
=



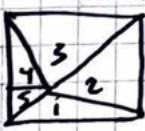
↳ same triangulation

Geometrically Different

Geometrically equivalent



=



Topologically Different

↳ different triangulation but same rectangle.

### Manifolds

\* Triangle is called manifold if:

mesh

(i) Intersections of two triangles are (i.e. non-empty).

~~empty set~~

~~manifold~~

~~each corner~~

~~edge has at most one vertex~~

~~Non-manifold~~

→ manifold (i.e. no empty intersection) & common vertex

→ manifold (i.e. common edge).

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

vertex is connected to two closed fans



closed fans = Non-manifold

(iii) Vertices:

- \* Build a single open fan → border vertex
- \* Build a single closed fan → inner vertex

vertex connected to an open fan

→ non-manifold

### Mesh Data Structures

(i) Shared Vertex

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

Shared  
Vertex  
data  
structure

+ List v → List f.

v  
vn  
f

--> 3 vertices

--> 3 faces

\*(ii) Half-Edge Data Structure

no only accepts manifold representations (since stores adjacent edges)

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

Useful for geometric computations

BRUNNEN

## Surfaces

for terrain in games

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

Constructed from Bezier Curves  $\rightarrow$  Bezier patches

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

Bezier curve

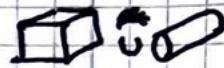


Bezier patch:

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

- (iii) Constructive Solid Geometry:

Boundaries created from boolean operations:



- (iv) Implicit Surfaces:  $\rightarrow$  SDF  
 $\rightarrow$  Heat  
~~mass~~  
 $\rightarrow$  Density

## Defining Implicit Surfaces:

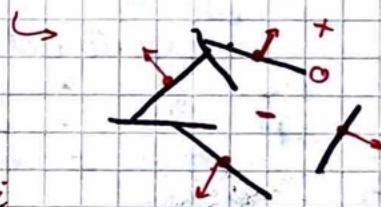
- local method
- bad with outliers
- global method
- better at handling outliers

(i) Hoppe method:  $f(\vec{x}) = (\vec{x} - \vec{p}) \cdot \vec{n}_p \rightarrow$  normal of point  $\vec{x}$   
for any  $\vec{x} \in$   $\vec{p}$   $\rightarrow$  nearest  $p$ . (from point cloud e.g.)

$\rightarrow$  creates Voronoi diagram which defines closest points

$\rightarrow$  Piecewise linear

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

$$\left. \begin{array}{l} f(\vec{p}_i + \epsilon \vec{n}_p) = \epsilon \\ f(\vec{p}_i - \epsilon \vec{n}_p) = -\epsilon \end{array} \right\} \text{our constraints}$$

some C2 function

$$\left. \begin{array}{l} \alpha_1 | \alpha_2 | \dots | \alpha_n | b | d \\ n+3+1 \end{array} \right\} \text{Unknowns: } n+4$$

Linear  $\Rightarrow$  SVD to solve for  $\alpha$

## (iii) Poisson Reconstruction

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

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.

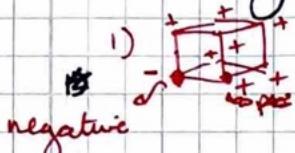


- 2) Sphere Tracing:
  - \* Dynamic Steps
  - \* Stop if distance below threshold.

(Also-to method usually).



- 3) Marching Cubes:
  - 1) Determine zero crossings for all grid cells.



- 2) If zero-crossing use Lookup table to triangulate

- 3) Linearly interpolate triangulation

- 4) Use which vertex negative as ID in Lookup Table.

- 5) ~~Find~~ Find Triangulation

- 6) Linearly interpolate triangulation

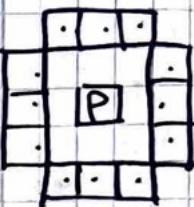
## Lecture 3: 3D reconstruction Methods

### Feature Detection

(i) Fast Detectors: for pixel  $p(x, y)$  if  $\geq 12$  pixels are brighter

⇒ feature points      partial derivatives w.r.t.  $x, y: I(u, v) - I(u+1, v)$   
 $\frac{\partial I}{\partial x}$

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



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

we find  $\lambda_1, \lambda_2$  of  $H$ :  $\lambda_1 \sim \lambda_2 + \text{small} \Rightarrow \text{flat region}$   
 $\lambda_1 \gg \lambda_2 \Rightarrow \text{OR } \lambda_2 \gg \lambda_1 \Rightarrow \text{Edge region}$   
 $\lambda_1 \sim \lambda_2 + \text{large} \Rightarrow \text{corner regions.}$

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

if $\lambda_1 \sim \lambda_2 \sim 0 \Rightarrow R \sim 0$ ⇒ i.e. flat region	or	if: $ R $ small ⇒ flat region
if $\lambda_1 \gg \lambda_2 \Rightarrow \lambda_1 \lambda_2 \sim 0 \Rightarrow R = -\alpha \lambda_1^2 \Rightarrow$ edge	or	$R \ll 0 \Rightarrow$ edge regions
if $\lambda_1 \sim \lambda_2 \gg 0 \Rightarrow \lambda_1 \lambda_2 - \alpha \lambda_1^2 \Rightarrow$ corner.	or	$R \gg 0 \Rightarrow$ corner regions

↳ usually  $0.04-0.08$

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

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

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

(iii) ConvNet: Basically siamese (shared weights) networks  
 inputting some image/volume patches 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:

→ 2 cameras

→ we know each camera's intrinsics

→ we know each camera's extrinsics

→ we segment subject in each camera

Carve over voxel grid based on  
 each camera segmentation

(ii) Structure from Motion (SfM): Given: → Set of images  
 → Features + correspondences  
 for each image

### Bundle Adjustment

Re-projection error:

$$\text{optimizes over } E_{\text{re-proj}}(T, X) =$$

$$T_i^* + \sum_{j=1}^m \sum_{i=1}^n \|x_{ij} - \pi_i(T_i X_j)\|_2^2$$

$X_j$

Images  $i \rightarrow$

feature points  $j$  of feature point  $j$  in image  $i$   $\rightarrow$  2D world  $\rightarrow$  3D  $X_j$  transformed to camera frame  $i$ :  $T_i X_j$  then projected to camera  $i$

Bundle Adjustment

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

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

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

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

(iii) Multiple View Stereo (MVS): Given: → Set of images

→ Poses of ~~all~~ images  $T_i$   
 → 3D locations of features  $X_j$

Goal: → Dense Point Cloud

- \* We can do frame pair matching + vote

- \* Global optimization is NP-Hard

- \* Plug point clouds to surface fitting e.g. poisson surface fitting

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

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

Frame to  $m$

Frame OR

Frame to Model

- \* Basically real-time SfM + MVS ~~after~~

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

- \* Types:
  - 1) ORB SLAM
  - 2) LSD SLAM.

} can also either use feature points for alignment or directly input

RGB-D { 3) KinectFusion

SLAM usually has Loop-Closure Processing

RGB-D Bundling

$$E_{\text{bundle}}(T) = \sum_{i,j} \sum_n \|T_i p_{ik} - T_j p_{jk}\|_2^2$$

$$E_{\text{depth}}(T) = \sum_{i,j} \sum_n \|(p_k - T_i^* T_j^* \pi_j(D_j(T_j(T_i p_k)))) \cdot n_n\|_2$$

$$E_{\text{color}}(T) = \sum_{i,j} \sum_n \|\nabla I(T_i(p_k)) - \nabla I(\pi_j(T_j^{-1} T_i p_k))\|_2$$

Color gradient of original pixel  $p_k$  ↪ color gradient of projected pixel  $\pi_j$

Take back-projected & corr. in  $i$ ,

2) Transform to  $g_j$  w/  $T_j$

3) to  $j$  camera frame

4) Project to 2D

5) Depth Value  $n_j$

6) Back-project  $\pi_j$

7) Transform to  $g_j$  w/  $T_j$

8) to original image frame  $T_i$

## Lecture 4: Optimization Methods:

### Linear Least Squares

→ Means the model is linear in parameters

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

where  $X_2$  can be  $x^2$ ,  $\log x$ ,  $\frac{1}{x}$  doesn't matter  
still linear since linearly multiplying by  $\beta_2$ .

\* 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.  $Y = \beta_0 X^{\beta_1}$   $\beta_1$  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

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

(ii) Gauss-Newton: Approximation of Newton:

$$H_f \approx 2 J_F^T J_F \Rightarrow x_{k+1} = x_k - [2 J_F^T J_F]^{-1} \nabla f(x_k)$$

$$\Rightarrow f(x+t) = \nabla f + H_f \cdot t = 0$$

assume  $x_k \in \mathbb{R}^n$   $J_F \in \mathbb{R}^{n \times n}$  residuals where  $F \in \mathbb{R}^m$

$$\Rightarrow t = H_f^{-1} \nabla f$$

$$\text{so: } x_{k+1} = x_k - t.$$

$$\underbrace{2 J_F^T J_F}_{A} \underbrace{[x_k - x_{k+1}]}_x = \underbrace{\nabla f(x_k)}_b$$

↳ (CPG)

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

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

↳ Tikhonov regularization

i.e. damping factor for each term in  $J$  s.t.  $\lambda \rightarrow 0$  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  $\lambda$  scaled by curvature of (diagonal is independent parts of curve)

(v) BFGS / LBFGS:

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

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

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:

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

(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:  $\frac{df(x)}{dx} \approx \frac{f(x+h) - f(x)}{h}$   $h \rightarrow 0$  i.e. take very small  $h$ .

$\rightarrow$  Easy to implement  $\Rightarrow$  good for debugging

$\rightarrow$  Slow and numerically unstable

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

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

(iii) Symbolic Differentiation:

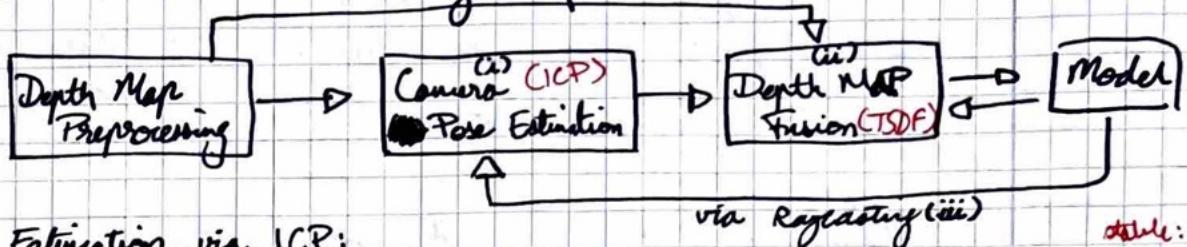
Basically no numerical solutions (Maple :))

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

2) Gradient-Free: Monte Carlo etc.

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

## Kinect Fusion Reminder weight updates



### (i) Pose Estimation via ICP:

Basic Steps of ICP:

if coarse  
to fine registration  
for conciliation  
problem.  
iterate

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

Error:  

- Use all points
- Uniform subsampling
- Random sampling
- Stable sampling



### Finding Correspondences in ICP:

only works  
if we have  
depth values

- a) Nearest neighbour → FLANN etc.
- b) Projective Correspondences → did in da project i.e.

\* We can improve matching via  
only checking compatible points:

$\pi(T \cdot \pi^{-1}(p_2))$   
my correspondence  
in other frame.

- e.g. at border.
- a) Color Compatibility
- b) Normal Compatibility
- c) Pruning
- d) Reject angle and distance outliers

### Errors in ICP:

has closed  
form linear  
solution !!

- a) Point to Point:  $\sum_i \|T p_i - q_i\|_2^2 = E(T)$
- b) Point to Plane:  $\sum_i \|(T p_i - q_i) \cdot n_i\|_2^2 = E(T)$
- \* Usually solved w/ LM if not linearized
- \* Symmetry ICP:  $C_p - q = (n_p + n_q)$

idea here  
is energy  
lowest if  
 $p_i \in q_i$  or  
the same surface  
w/ normal  
 $n_i$

↳ can be linearized in  $\approx 0$   
assuming small rotations  $\cos \theta \approx 1$   
→ Cholesky / SVD

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

- Basically representation of distance to zero-crossing band.
- TSDF: just applies truncation, i.e. if  $> 1$  from zero-crossing
- /  $< -1 \Rightarrow +1 / -1$  everywhere

Volumetric Fusion: Surfel representation: point + normal + radius  
also exists current depth based on our current depth value

weights of running average concept:  $D_{i+1}(x) = w_i(x) \cdot D_i(x) + w_{i+1}(x) \cdot d_{i+1}(x)$   
 can either be 1  
 everywhere OR lower weight, away weights currently (i.e. distance to iso surface)  $w_{i+1} = w_i(x) + w_{i+1}(x)$   
 to camera + weight, away weights (uncertainty weighted)

### (iii) Surface Rendering

Important to note on SDFs:

=0 → surface

<0 → unknown

(would have another surface ~~there~~ behind)

>0 → free space.

a) Raycasting: Throw multiple rays from camera.

→ Surface found from sign change + interpolate to find intersection

→ Normal found from VTSDR.

b) posteriorizing

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 reintegrate previous 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

- ~~slow~~ lookup slower = raycasting not as fast.

+ cheap streaming

### Loop Closure Follow Ups:

#### 1) Elastic Fusion:

→ Surface based

→ Detects loop closure by localization

→ Non-rigidly warp to close loops.

#### 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

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

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

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

minimizes distance between vertex and target

1-ring fan:



### Regularizer Types

→ Our goal for regularizers: → smooth global shape  
→ Global shapes preserve local detail  
→ Should be zero in undeformed state.

#### i) Laplacian Surface Editing:

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

→ our fan neighbourhood (1-ring)

→ i.e. deformation of  $v_j'$  should be very similar to 1-ring neighbours  $v_i'$

#### ii) As Rigid as Possible (ARAP):

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

done per fan over all vertices.

means that after some  $R$  ~~is~~ edge

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

$$v_i \quad v_i'$$

$$v_j \quad v_j'$$

after some rotation.

\* we could add weight to give bigger priority to closer neighbours:

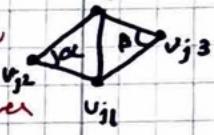
$$w_{ij} = \frac{1}{2} (\cot \alpha_{ij} + \cot \beta_{ij})$$

3 positions per vertex → 3 Euler angles

as  $\cot \alpha_{ij} + \cot \beta_{ij}$

→ if  $v_{ij}$  closer

then  $v_{ij}'$  means lower



\* Solving: Unknowns:  $3n + 3n$  non-vertices

(per position)

Constraints:  $3n$  valence +  $3c$

weight for  $v_{ij}'$

average a)  
# for neighbours

↳ from  $E_{fit}$

→ 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.

(ED)

iii) Embedded Deformation: relaxed version of ~~ARAP~~ ARAP

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

$$Err_{rot} = \sum_i Rot(M_i), \quad Rot(M) = (C_1 \cdot C_2)^2 + (C_1 \cdot C_3)^2 + \dots + (C_1 \cdot C_{i-1})^2 + \dots$$

\* 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 Phases → instead of deforming mesh = computationally expensive deform proxy

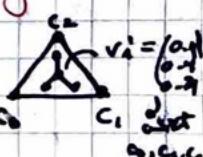
we can apply →

i) Cage: linear combination of vertices w.r.t. cage.

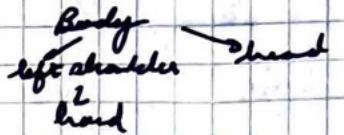
(we use barycentric coordinates for that.)

ARAP/ED  
deformation here.

BRUNNEN → we deform cage and by proxy our vertices using barycentric coordinates.



ii) 3D Grids: ~~discrete~~

e.g. Blend or iii) Skeletons: hierarchical: 

Called  
linear  
blend skinning

{  
- Each vertex where there is a linear combination of bone  
matrices:  $v_i = \sum w_{i,j} M_j v_i'$   
weight for each bone  $\rightarrow$  each bone has a matrix  
each bone

iv) Deformation Graphs: assigns a region/neighborhood a central  
node.  $\rightarrow$  similar to cage but instead we have deformation  
zone.

Non-Rigid Tracking  $\rightarrow$  find correspondences then ICP.

needs predefined  
mesh.

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

$$E = \underbrace{\lambda_{\text{point}} E_{\text{point}} + \lambda_{\text{plane}} E_{\text{plane}} + \lambda_{\text{color}} E_{\text{color}} + \lambda_{\text{reg}} E_{\text{reg}}}_{\text{from regular point to plane ICP}}$$

i.e. lower energy if colors match

$$E_{\text{point}}: \sum_i \|v_i - T v_i'\|_2^2$$

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

$$E_{\text{color}}: \sum_i [I(\pi(v_i)) - I'(\pi(T v_i))]^2$$

$$E_{\text{reg}}: \text{ARAP.} \quad \text{lower energy if colors match.}$$

\* Important hyperparameter-thresholds for non-rigid ICP:

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

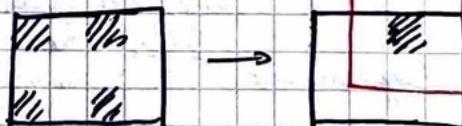
Happens with  
more cameras.  $\rightarrow$

??

\* Correspondence Association:

$\rightarrow$  for speed.

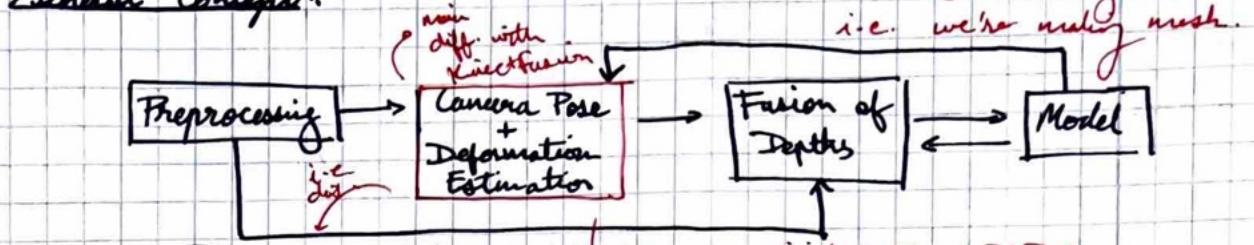
- c.g.:  
1) Projective correspondences every k-th element  
2) Kernel around k-th sample for minimum reduction



no 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 losses  
course is

Main Methods discussed:

th reconstruction: (i) DynamiC Fusion: ↳ Uses Deformation Graphs (ED)  
course is ↳ Non-rigid ICP

↳ Volumetric Fusion.

than this

surface moving ↳ More stable  
skin deforming? ↳ Better texture  
| (texture) We can increase resolution based on grid  
| This means higher stability compared to dynamic fusion

(ii) Volume Deform:

- ↳ uses regular deformation grids + ARAP
- ↳ Non-rigid ICP using sparse SIFT matches.
- ↳ Volumetric Fusion

More stable  
+ Better texture

(iii) Real-time Geometry & Motion Reconstruction:

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

(iv) 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.

(v) Motion2Fusion:

- ↳ Basically faster Fusion 4D: 200FPS

↳ ED Graphs

↳ Key Volumes

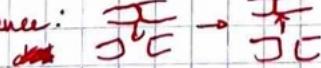
↳ Multiple Depth Cameras.

↳ 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 data to reference:



2) Detail layer: they increased resolution of grid at surface

close to open: open to close:

reference to data to data to reference:

(1) (2)

(viewing rays)

(vi) Lookin Good: ↳ increase / improve reconstructions using deep architecture

(vii) DeepDeform: ↳ find better correspondences than SIFT for best tracking

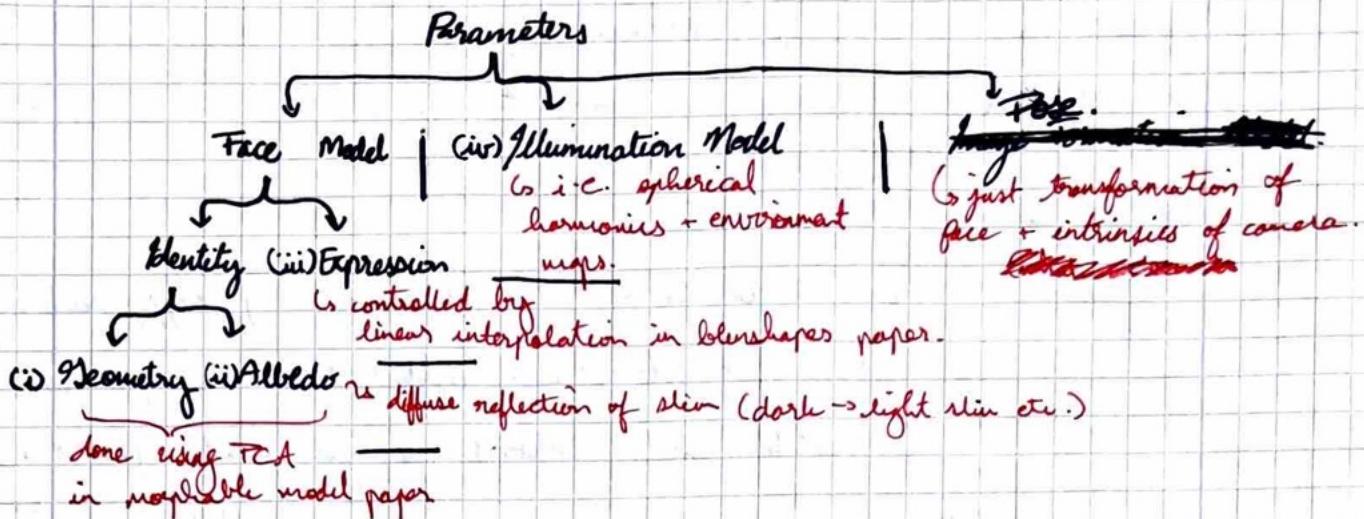
The two RGBs ↳ ↳ use siamese network that shares weights for two  
are the source ↳ RGBs to give you correspondences.

query points and  
target where correspondences  
need to be found.

$\rightarrow \Rightarrow$  For non-rigid reconstruction: no 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 n~~: 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:

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

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 has no semantic meaning though /  
captures data. 2) BlendShapes <sup>has semantic meaning</sup>

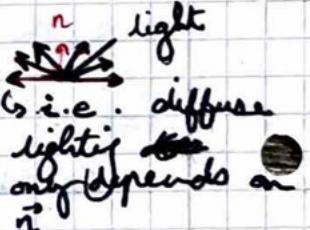
- \* BlendShapes: Steps:

- 1) Artist makes n number of target expressions (e.g. happy and 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  $\rightsquigarrow$  parameterization of environment maps
- \* Environment maps place the object a cube/sphere map of source panorama scene
- 2) Project lighting onto object (which is in center of sphere/cube)
  - $\rightsquigarrow$  Assumptions:  $\rightarrow$  Distant light  
 $\rightarrow$  No scattering
- \* Spherical harmonics parametrize/reduce dimensionality of environment maps. (orthogonal basis)
  - $\Rightarrow$  we can linearly interpolate basis spheres for different lighting
  - $\rightsquigarrow$  Assumptions:  $\rightarrow$  Lambertian Surface:  
 $\rightarrow$  Distant Smooth lighting



## Lecture 9: Face Tracking and Reconstruction

\* Goal: find parameters  $\hat{P}$  of parametric face model to fit RGB 2D image.

\* We make energy term:  $\rightarrow$  for energy minimization

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

Energy of color and geometry of mesh wrt.  $P$       Correspondence terms wrt.  $P$       regularization of  $P$

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

point to point       $\rightarrow$  colors.

$$E_{geom}(P) = \sum_i \| M_i(P) - D(\pi(M_i(P))) \|_2^2 + [(M_i(P) - D(\pi(M_i(P)))) \cdot n_i]^2$$

point to plane      point to point

$$E_{geom}(P) = \sum_i \| M_i(P) - D(\pi(M_i(P))) \|_2^2 + [ (M_i(P) - D(\pi(M_i(P)))) \cdot n_i ]^2$$

summed over all vertices      depth vertex      correspondant depth value at projected mesh vertex in camera frame      dot product of mesh mesh normal.

$$E_{col}(P) = \sum_i \| M_i(P) - I(\pi(M_i(P))) \|_2^2$$

colors of mesh at vertex  $i$  (in mesh)      colors of pixel at projected mesh vertex in camera frame.

$$E_{reg}(P) = \sum_k \left( \frac{\alpha}{\sigma_{id,k}} \right)^2 + \sum_l \left( \frac{\beta}{\sigma_{all,l}} \right)^2 + \sum_m \left( \frac{\gamma}{\sigma_{exp,m}} \right)^2$$

done for all correspondences

$$E_{sparse}(P) = \sum_j \| \pi(M_j(P))_m - C_j \|_2^2$$

2D coordinates of projected mesh in camera frame       $\rightarrow$  2D coordinates of correspondences in 2D image

$\alpha, \beta, \gamma$  are parameter vectors of id, all, exp & centered at zero of. The summands done for each dimension

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

$\rightarrow$  Gauss Newton or LM solver.

$$\Rightarrow P_{n+1} = P_n - (\mathbf{J}_F^\top \mathbf{J}_F)^{-1} \mathbf{J}_F^\top \cdot F \quad \text{as Usually done Rough} \rightarrow \cancel{\text{Pyramid Dense}}$$

Partial Derivatives of 2D image wrt.  $P$ :

$\rightarrow$  needed for  $E_{color}$  +  $E_{shape}$

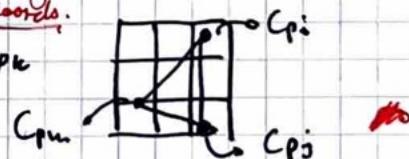
colors at  $p_i$  etc.

1) For each color pixel:

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

Barycentric coords.

linear interpolation of vertex colors.



$\rightarrow$  probably the same but not sure :/

2) record triangle  $v$ , barycentric coords  $\alpha, \beta, \gamma$ .

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

$\rightarrow$  This parameter update done:

- \* Can be done without depth via multiple views
  - ↳ energy jointly optimized for 6 images
- \* Energy optimization dependent on:

- Input Data (R<sub>b</sub>B vs. R<sub>a</sub>B-D)
- Face Model + Illumination (coarse vs. dense)
- Run-time requirements, hardware
- Application

- \* Two main papers using dso:
  - Real-time expression transfer for facial reenactment
    - ↳ for two live videos
  - Face2Face
    - ↳ 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

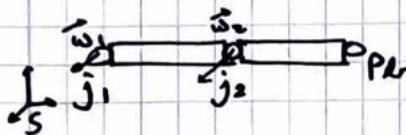
2) PCA of all mesh vectors.

SCAPE: → joint shape + pose learning.

→ blend shapes used first for poses.



### Skeleton Siderote:



each joint has a transformation matrix related to rotations by Rodrigues formula

$$\Rightarrow p_s = g_s(\vec{w}_1, j_1) g_s(\vec{w}_2, j_2) p_b$$

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

Linear Blend Shading: For any vertex  $t_i$  on mesh:

vertices are linear combinations of all joints  $\in$   $t_i' = \sum_{k=1}^n w_{k,i} g_k(\theta, \vec{j}) t_i$

weight given to transformation of each joint

Corrective Blend Shapes: each joint

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

Non-rigid registration:  $E = \underbrace{w_d E_d}_{E_d = \sum w_i \text{dist}^2(T_i, v_i, D)}$  +  $\underbrace{w_s E_s}_{E_s = \sum \|T_i - T_j\|_F^2}$  +  $\underbrace{w_m E_m}_{E_m = \sum \|T_i v_i - m\|^2}$

Data error.

Transformation Smoothness

marker Errors.

### Hand Models:

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

$\rightarrow C_2$  continuity

$\rightarrow$  Fast convergence

## Issues with Model Generation + Fitting:-

- Experience Scanning setup.
- Neutral Pose similar but not identical
- Ambiguities in fitting  $\curvearrowright$  change shape or expression?
- Drift in non-rigid registration