

## Concepts

### 3D Representations:

Voxels: binary (or probabilistic) occupancy grid. Cubic memory

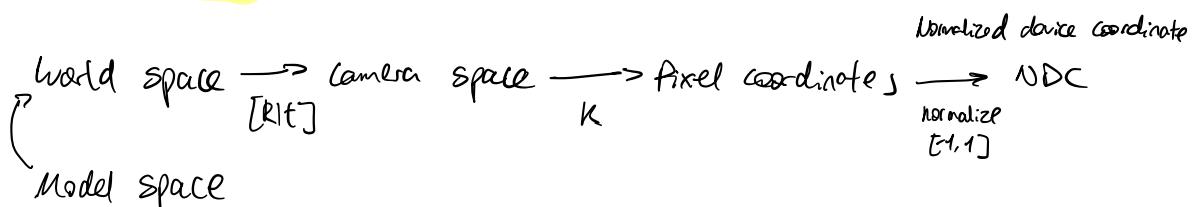
Point clouds: 3D points (unstructured) with attributes

Polygon Mesh: set of vertices and faces (color = interpolation of face vertices' color)

Parametric surface: set of control points (splines)

Implicit surface: signed distance field. Implicit:  $x+y=0$   
Explicit:  $f(x) = -x$

## Camera Spaces



## Extrinsics

$$E = \begin{bmatrix} R & t \\ 0 & 0 & 1 \end{bmatrix} \quad R \in \mathbb{R}^{3 \times 3} \in SO(3) \quad t \in \mathbb{R}^3$$

$$\begin{array}{l} \rightarrow |R|=1 \\ \rightarrow RR^T=I \end{array}$$

6 DoFs: 3 for  $R$  (in  $SO(3)$  lie algebra representation)  
3 for  $t$

$$R = R_z(\gamma) R_y(\beta) R_x(\alpha)$$

right up look position  
↓ ↓ ↓  
 $\text{Camera Pose} = \begin{pmatrix} R_x & U_x & L_x & X_c \\ R_y & U_y & L_y & Y_c \\ R_z & U_z & L_z & Z_c \\ 0 & 0 & 0 & 1 \end{pmatrix}$

$E = \text{Pose}^{-1}$ . This definition is relative. Could also be  $E = \text{Pose}$

## Intrinsics

$$K = \begin{pmatrix} f_x & r & m_x \\ 0 & f_y & m_y \\ 0 & 0 & 1 \end{pmatrix}$$

$$f_x = f \frac{w}{W} \quad f_y = f \frac{h}{H}$$

$f$  = focal length in mm  
 $w, h$  = pixels

$W, H$  = sensor width/length in mm

$f$ : focal length, dist from pinhole to image plane,  
aspect ratio:  $f_x/f_y$ . In pixel units

$m$ : principal point, dist to proj. center ( $w/2, h/2$ ).  
In pixel units

$r$ : axis skew

## Calibration

Radial Distortion: fish-eye effect.



$$x_{\text{corrected}} = x(1 + k_1 r^2 + k_2 r^4 + k_3 r^6)$$
$$y_{\text{corrected}} = y(1 + k_1 r^2 + k_2 r^4 + k_3 r^6)$$

Tangential Distortion: barrel effect



$$x_{\text{corrected}} = x + [2p_1 xy + p_2(r^2 + 2x^2)]$$
$$y_{\text{corrected}} = y + [p_1(r^2 + 2y^2) + 2p_2xy]$$

## Stereo Matching

Estimate depth with 2 rgb cameras.

- Rectification: make epipolar lines parallel to horizontal axis of image planes.
- Search: find match (with 2D features) along epipolar lines

Epipolar constraints,  
essential matrix.

$$SSD = \sum_{(i,j)} [I_{\text{left}}(i,j) - I_{\text{right}}(i,j)]^2$$

Possible features: SIFT, SURF...

Window size: small  $\rightarrow$  noisy  
large  $\rightarrow$  few detail

Search line also correlated to window size

Baseline (distance between cams): small  $\rightarrow$  less accurate (small disparity results in large error)  
large  $\rightarrow$  hard to find matches (occlusion)

## Devices

Passive:

Stereo (Multi-) camera: rely on 2D features and expensive matching / finding

Active:

TOF: IR spectrum (fails outdoor), sensitive to scattering, indirect light

Structured light: IR spectrum (fails outdoor), need calibration between projector and sensor

Laser scanner / Lider: slow and precise or fast and sparse.  
precision due to trivial feature matching.

## Surface Representations

### Polygonal Mesh

Approximation error  $O(n^2)$

Arbitrary topology (triangulation)

Efficient rendering (GPU rasterization)

If it's a manifold if:

- Any intersection of 2 triangles
  - Empty
  - Common edge
  - Common vertex

- Edges have either 1 or 2 adjacent triangles

- For a vertex the adjacent triangles build a single open or closed fan topology (triangulation), geometry (shape)

Formats: obj, off, ply ...

Half-edge data structure: easy geometric query (directed edges)

### Explicit Surface

$f(x,y) : \mathbb{R}^2 \rightarrow \mathbb{R}$

$\delta(x,y) = (x,y, f(x,y))$

only one height value per  $(x,y)$  pair  
e.g. height map

### Parametric Surface

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

freeform surface with set of control points  
e.g. Bezier, splines

### Constructive Geometry

surface created by boolean operations on primitive solids. OP: U, ∩, -, ...

## Implicit Surface

$$f(x, y, z) : \mathbb{R}^3 \rightarrow \mathbb{R} \quad \text{SDF (signed distance function)}$$

$$\text{Hesse normal form: } f(x, y, z) = \left( \begin{pmatrix} x \\ y \\ z \end{pmatrix} - \vec{p} \right) \vec{n} = 0$$

Given sample points find SDF?

### Hoppe

$$f(\vec{x}) = (\vec{x} - \vec{p}) \cdot \vec{n}_p \quad \vec{p} \text{ the nearest neighbor of } \vec{x} \text{ in the given sample points}$$

↳  $f(\vec{x})$ : Distance to the tangent plane of  $\vec{p}$

piecewise linear (linear boundaries, Voronoi diagram), not smooth.  
Can't handle noise, outliers...

### KBF

theo: complex functions can be approximated as the sum of simple scaled and translated kernel functions  $\varphi(x)$

$$f(\vec{x}) = \sum_i d_i \cdot \varphi_i(\vec{x}^3) + \vec{b} \cdot \vec{x} + d \quad \varphi_i(\vec{x}^3) = \|\vec{p}_i - \vec{x}\|^3 \quad \text{biharmonic basis function}$$

i: sample points

params (unknowns):  $\begin{matrix} d_i & (n) \\ \vec{b} & (3) \\ d & (1) \end{matrix} \quad \left. \begin{matrix} \\ \\ \end{matrix} \right\} n+4$

constraints:  $\begin{matrix} f(\vec{p}_i) = 0 & (n) \\ f(\vec{q}_i) = \pm \epsilon & (n) \end{matrix} \quad \left. \begin{matrix} \\ \\ \end{matrix} \right\} 2n$

$$\vec{q}_i = \vec{p}_i \pm \epsilon \cdot \vec{n}_{\vec{p}_i}$$

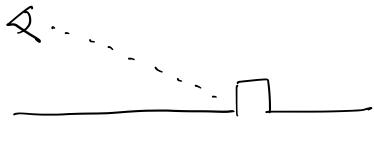
$$A \cdot \vec{x} = \vec{b} \quad \text{overdetermined} \\ (A^T A \cdot \vec{x} = A^T \vec{b})$$

$$\begin{bmatrix} \vdots & \vdots & \vdots \\ - & - & - \\ \vec{p}_i & \vec{q}_i & \vec{p}_i \\ \vdots & \vdots & \vdots \\ \vec{b} & \vec{b} & \vec{b} \\ \vdots & \vdots & \vdots \\ d & d & d \end{bmatrix} \begin{bmatrix} d_1 \\ \vdots \\ d_n \\ b_1 \\ \vdots \\ b_3 \\ d \end{bmatrix} = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ \pm \epsilon \\ \vdots \\ \pm \epsilon \\ \vdots \end{bmatrix}$$

Can handle noise and outliers, but only works with small point set (slow).  
A workaround is to partition point cloud into parts and fit individually.  
Or coarse-to-fine ...

## Rendering SDF

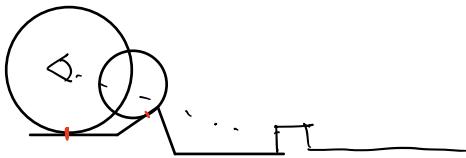
### Ray Marching



fixed steps, linear interpolation when crossing occurs.

march forward to find intersection

### Sphere Tracing



A type of ray marching for SDFs.  
At each step solve for  $f(x,y,z) = 0$  (max sphere radius).

Then we can take step of at least this value (safely).  $\rightarrow$  dynamic steps

stop if distance is below some threshold.

### Marching Cubes

Create a polygonal mesh from a SDF. In 3D each SDF cell has 8 vertices  $\Rightarrow$  check zero crossing in lookup table of  $2^8 = 256$  combinations.

vertex position can be adjusted with linear interpolation of SDF values.

## Poisson Surface Reconstruction

Goal: obtain mesh (extracted from isosurface) from input sample points and normals.  $O(n^2)$  in time and memory. Need dense point cloud.

### Estimate indicator function $X$

$$\nabla X = \vec{v} \Rightarrow \Delta X = \nabla \cdot \nabla X = \nabla \cdot \vec{v} \quad (\text{Laplacian} = \text{divergence of gradient})$$

vector field  $\vec{v}$  is the surface integral (we don't know it before hand), but it can be approximated by the summation over patches (neighborhood).

$$\vec{v}(q) \approx \sum_{s \in S} |P_s| \cdot F_{sp}(q^s) \cdot s \cdot \vec{N} \quad \vec{v}(q) = \nabla(X * f)(q) \quad f = \text{Gaussian filter}$$

$$\text{Solve for } \sum_{\substack{o \in \mathcal{O} \\ L \geq \text{dim}}} \| \langle \Delta X, f_o \rangle - \langle \nabla \cdot \vec{v}, f_o \rangle \|^2 \Rightarrow \min_{X \in \mathbb{R}^{|\mathcal{O}|}} \| Lx - v \|^2 \quad \text{Gauss-Seidel}$$

### Mesh reconstruction from $\chi_m$

Evaluate  $\chi$  at sample points and take average to find appropriate isovalue.  
Perform marching cubes to extract surface (mesh).

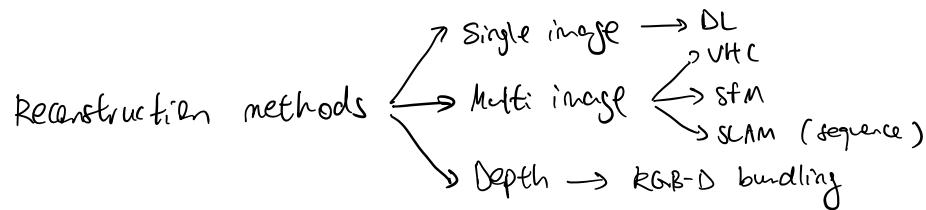
### Screened Poisson

$\mathcal{O}(|P| \log |P|) \approx \mathcal{O}(n^{1.5})$  complexity

$$(\Delta - \alpha \tilde{\mathbf{I}}) \chi = \nabla \cdot \vec{V}$$

$\alpha$  = screening weight

## Reconstruction Methods



### Single Image

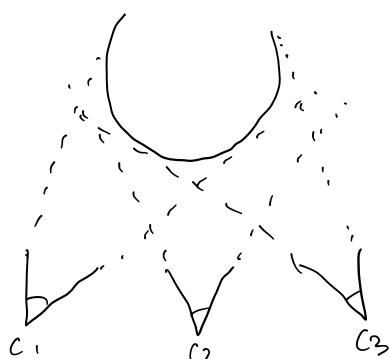
Extremely ill-posed problem, often need strong priors, guess shape from shading.

Blurry shape  $\rightarrow$  deep learning

### Multiple Image

#### Visual Hull Craving

Calibrated multi-camera setup (known  $E, K$ )



Need object segmentation to identify foreground (silhouette). Then with known  $(E, K)$  back-project object to 3D as a silhouette cone. Intersect the cones from other views.

Apply smoothing kernel to improve quality.

Could also do with depth data by fusing views

Sensitive to noise

#### Structure from motion (SfM)

Construct mesh from images taken from different viewpoints.

1. Estimate camera params  $(E, K)$  and 3D points.

1.1 Find keypoints of each image

- PAST (features from accelerated segment test): compare the intensity of  $p + \theta$  to the 16 nearby pixels forming the circle.

If more than  $n$  (usually 12) contiguous pixels are all greater or all smaller than  $I_{p+th}$  then  $p$  is a corner (Keypoint).

- HARRIS corner detector: find corners in a  $3 \times 3$  sliding window manner.  $E(u,v)$  measures the sum of SSD in all 8 directions.

$$E(u,v) = \sum_{x,y} w(x,y) [I(x+u, y+v) - I(x,y)] \approx [u, v] M [u, v]$$

$$M = \sum_{x,y} w(x,y) \begin{bmatrix} I_x^2 & 2I_x I_y \\ 2I_x I_y & I_y^2 \end{bmatrix}$$

→ eigenvectors give direction of greatest change in SSD  
 ↳ eigenvectors  $(\lambda_1, \lambda_2)$  can be used to determine:

- $\lambda_1, \lambda_2$  small  $\Rightarrow$  flat region
- $\lambda_1, \lambda_2$  large  $\Rightarrow$  corner
- $\lambda_1, \lambda_2$  differs  $\Rightarrow$  edge

## 1.2. Correspondence matching (feature descriptor)

Naively SSD. But ideally should be: scale, view, light invariant

- SIFT (scale invariant feature transform)

Detector:

Scale invariance  $\Rightarrow$  different octaves each one's size is half of prev.  
 DoG  $\Rightarrow$  progressive blur images (change  $\sigma$ ) in a octave and compute difference of gaussian (DoG)

Laplacian of gaussian  $\Rightarrow$  a pixel is compared to its neighbor (8), 9 from prev and 9 from next scale.  
 If this pixel is maxima, then it's a potential keypoint.

Descriptor:

Scale, rotation, illumination invariant.

$16 \times 16$  window around the keypoint which is then subdivided in  $4 \times 4$  sub-blocks, each one with 8-bin orientation histogram  $\Rightarrow 4 \times 4 \times 8 = 128$  dim.  
 Matching  $\Rightarrow$  dot prod.

rot invariance: from all orientation extract keypoint's orientation

illu invariance: cap (clamp) the value in feature vector and normalize.

- ConvNet (pretrained weights, ReNet)

### 1.3 Bundle Adjustment

formulate the reprojection errors and minimize it by solving the non-linear least square problem.

$m$  images (cams),  $n$  total keypoints :  $6(m-1)$  extrinsics  
 $3n$  points  
make sure  $2mn > 6(m-1) + 3n + (\text{intrinsic})$  possibly intrinsic } params

solve to obtain sparse point cloud and camera pos.  $2m \cdot n$  (2D) } constraints

$$E_{\text{proj}}(T, X) = \sum_i^m \sum_j^n \|x_{ij} - T_i(T_i \cdot x_j)\|_2^2$$

### 2. Multiview Stereo (MVS)

With result from 1. ( $K, E$ , points), estimate dense point cloud.  
Do dense stereo matching between frame pairs (note for results)

### 3. Reconstruction

With dense point cloud from 2. estimate normals with PCA and do Poisson reconstruction to obtain mesh. Texture resulting surface

PCA for normals: with the neighbors of a point compute centered cov. matrix and perform eigen decomposition. Eigen vector corresponding to 3rd eigenvalue (ordered) is the surface normal

## Simultaneous Localization and Mapping (SLAM)

Similar idea to SfM+MVS but online.

Model is built incrementally by frame-to-frame tracking.

Bundle adjustment on keyframes only and use ORB features to save computing budget.

### RGB-D Bundling

Important: here  $T$ 's are camera-to-world matrix. (Extrinsics $^{-1}$ )

Depth data available, frame-to-frame tracking or frame-to-model (accumulate depth data into model, SDF). May have loop closure issue.

Sparse:  $E_{\text{bundle}}(T) = \sum_{i,j} \sum_k \|T_i p_{ik} - T_j p_{jk}\|_2^2$  ( $i, j$ ) frames,  $k$  correspondence,  $p_{ik}$  3D point (depth from sensor)

Dense:  $E_{\text{depth}}(T) = \sum_{i,j} \sum_k \|(\hat{p}_k - T_i^{-1} T_j T_{ik}^{-1} (D_j(T_{ik}(T_j^{-1} T_i p_k))) h_{ik})\|_2^2$   $k$  pixels,  $D_j$  depth map

$$E_{\text{color}}(T) = \sum_{i,j} \sum_k \|\nabla I(T_i p_k) - \nabla I(T_{ik}(T_j^{-1} T_i p_k))\|_2^2$$

## Optimization

### Least Squares

$$\text{Solves for } f(x) = \sum r_i(x)^2 \quad x^* = \arg \min_x f(x) = \arg \min_x \|F(x)\|_2^2 \quad f(x) = [r_1(x), \dots]^T$$

### Linear Least squares

residuals are linear w.r.t parameters. Could be expressed as

$$\begin{array}{c} Ax = b \\ \uparrow \quad \uparrow \\ \text{inpt} \quad \text{residuals} \end{array}$$

### Normal Equation

if  $Ax=b$  is overdetermined (#constraints > #params). solve by

$$ATAx = A^Tb \Rightarrow x = (A^TA)^{-1}A^Tb$$

### Iterative Methods

(parallelizable but slow convergence)

- Jacobi iteration :  $A = L + D + U$ ,  $x_{k+1} = D^{-1}(b - (L+U)x_k)$   
(sequential but fast convergence)
- Gauss-Seidel iteration :  $A = L + U$ ,  $x_{k+1} = L^{-1}(b - UX_k)$   
 $\left. \begin{array}{l} (L+U)x = b ; \\ Lx + Ux = b ; \\ Lx = b - UX ; \\ x^{k+1} = L^{-1}(b - UX^k) \end{array} \right\}$
- Conjugate gradient descent :  $x_{k+1} = x_k + \alpha_k p_k$      $\frac{1}{2}x^T Ax - b^T x$   
(PCG)                   $\alpha_k, p_k$  dependant to residual  $r_k$

### Direct methods

- QR :  $A = QR$ ,  $Q^T = Q^{-1}$   
 $\left. \begin{array}{l} PAx = Pb ; \quad Ly = Pb ; \\ y = L^{-1}Pb \Rightarrow x = U^{-1}L^{-1}Pb \end{array} \right\} \quad \begin{array}{l} Ax = y \\ PAx = Pb ; \quad Ly = Pb \Rightarrow Ly = Pb ; \\ y = L^{-1}Pb \Rightarrow x = U^{-1}L^{-1}Pb \end{array}$
- LU :  $PA = LU$
- Cholesky :  $A = LL^T$
- SVD :  $A = U\Sigma V^T$      $(V\Sigma V^T)x = b$ ;  $x = V\Sigma^{-1}U^T b$

### Non-Linear Least Squares

residuals non linear w.r.t params.

### 1<sup>st</sup> Order

Gradient Descent :  $x_{k+1} = x_k - t \nabla f(x_k)$

### 2<sup>nd</sup> Order

$$\text{Newton: } X_{k+1} = X_k - H_f(x_k)^{-1} \nabla f(x_k)$$

$$A \quad x \quad b$$

$$\text{Gauss-Newton: } H_f(x_k) \approx 2J_F^T J_F \quad X_{k+1} = X_k - (2J_F^T J_F)^{-1} \nabla f(x_k) \quad (2J_F^T J_F)(X_k - X_{k+1}) = \nabla f(x_k)$$

$$\text{Levenberg: } H_f(x_k) \approx 2J_F^T J_F + \lambda \cdot I$$

$$J_F = J_F(x_k)$$

$$\text{Levenberg-Marquardt: } H_f(x_k) \approx 2J_F^T J_F + \lambda \cdot \text{diag}(J_F^T J_F)$$

BFGS:  $H_f(x_k) \approx B_k$   $B_k(X_k - X_{k+1}) = \nabla f(x_k)$   $B_{k+1} = B_k - \alpha_k u u^T + \beta_k v v^T$ . In practice directly estimate  $B_k^{-1}$

L-BFGS  $\Rightarrow$  approximation with limited memory

## Handling Outliers

RANSAC: trial and error. for a pair of frames try different combination of correspondence (by only taking a subset) and see which gives the lowest error (outlier  $\Rightarrow$  high error)

Lifting Scheme: optimize additional parameter  $w$ , which are weights for each residual term. Let optimization decide which are outlier by assigning low value (ideally 0 for outlier, 1 for inliers). Regularization is added to avoid trivial solution.

$$f_{\text{rob}}(x, w) = \sum w_i^2 r_i(x)^2 + \lambda \text{reg} \sum (1-w_i)^2$$

Robust norms:  
(IRLS: Iterative Reweighted Least Squares)  
e.g.  $L_1$  is more robust to outliers than  $L_2$ , however hard to optimize. IRLS trick:

for any norm  $p$ :  $f(x) = \sum |r_i(x)|^p$  map to  $L_2$  in each iteration

$$\hookrightarrow \sum_{L_2} w_i r_i(x)^2, \quad w_i = |r_i(x)|^{p-2} \quad (p = p-2+2)$$

$\hookrightarrow$  changes at each iteration as  $x$  gets updated.

## Efficiency and Convergence

To make the optimization smoother and more convex, first apply smoothing and adopt coarse-to-fine strategy.

Exploit sparsity of  $J$  and parallelize optimization with GPU for maximal efficiency

## Derivatives

Jacobian :  $J_f(x) = \begin{bmatrix} \frac{\partial F}{\partial x_1} & \cdots \\ \vdots & \ddots \\ \frac{\partial F}{\partial x_n} \end{bmatrix}$

vars (n)

Gradient :  $\nabla f(x) = 2 (J_f(x)^T) F(x) \quad (\mathbb{R}^{m \times n})^T \mathbb{R}^{m \times 1} = \mathbb{R}^{n \times m} \mathbb{R}^{m \times 1} = \mathbb{R}^n$ .

When only 1 residual  $\Rightarrow J_f(x) = \nabla f(x)$ . Thus, gradient is a special case of jacobian

Hessian :  $H_f(x) = J_{\nabla f}(x)^T = \begin{bmatrix} \frac{\partial^2 F}{\partial x_1^2} & \cdots \\ \vdots & \ddots \\ \frac{\partial^2 F}{\partial x_n^2} \end{bmatrix}$

vars (n)

Ways to compute derivative:

## Numerical Derivative

$$\frac{df(x)}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

easy to implement and debug but slow and unstable

## Auto Diff

### Ceres Dual Numbers

choose infinitesimal unit  $e$ , such that  $e \neq 0$ ,  $e^2=0$ .

Let  $f(x)=x^2$ :

$$f(10+e) = (10+e)^2 = 100 + 20e + e^2 = 100 + 20e \Rightarrow f'(10) = 20.$$

Verification:  $f'(x)=2x$ ,  $f'(10)=2 \cdot 10 = 20$

## Symbolic Diff

Analyse compute graph at compile time, can simplify terms efficiently.  
However optimal solution is NP-complete

## Connection to DL

In DL mainly stochastic gradient descent is used to train the network.  
Done by backpropagation on compute graphs, often handled by frameworks like Pytorch. Stochasticity is need for large data set and also helps with local minima.

## Other Optimization methods

Inequality  $\begin{cases} \xrightarrow{\text{Lagrange Multiplier}} \\ \xrightarrow{\text{Primal Dual}} \end{cases}$

Differential eqns  
 Gradient-free  $\begin{cases} \xrightarrow{\text{Hence Carlo}} \\ \xrightarrow{\text{Hill Climbing, Simulated Annealing}} \\ \xrightarrow{\text{Genetic algorithm}} \end{cases}$

## Rigid Tracking and Reconstruction

Goal: RGB-D tracking, move camera over time and incrementally build the scene. Would need to align point clouds

### Aligning Point Clouds

#### Procrustes algorithm (Correspondence)

Given point clouds  $X, Y$ , want:  $\min_{R,t} \|Y - RX + t\|_F$

scaling: compute center for both point clouds. Scale objects to match the avg. distance from vertices to center of gravity in the target point cloud.

rotation:

1. Center  $X, Y \Rightarrow X_c = X - \bar{X}, Y_c = Y - \bar{Y}$   $\bar{x}, \bar{y}$  means
2.  $\min \|Y_c - RX_c\|_F^2 \Rightarrow \min \text{tr}(Y_c^T Y_c - Y_c^T X_c R^T - (X_c R^T)^T Y_c + (X_c R^T)^T (X_c R^T))$   
 $\|A\|_F^2 = \text{tr}(A^T A)$  doesn't affect  $R$   $R^T R = I$   
 $\Rightarrow \min \text{tr}(-2(Y_c^T X_c R^T)) \Rightarrow \max \text{tr}(U \Sigma V^T R^T) \Rightarrow \max \text{tr}(Z V^T R^T U)$   
 $Y_c^T X_c = U \Sigma V^T$  trace  
 $\Rightarrow \text{tr}(Z)$  already maximal, therefore find such  $R$  that cancels  $V^T$  and  $U \Rightarrow R = V U^T$  so  $\text{tr}(Z V^T R^T U) = \text{tr}(Z)$ .

But if  $|R| = |V U^T| = -1$  (mirroring)  $\Rightarrow R = V S U^T$   $S = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & -1 \end{pmatrix}$   
to flip sign

translation:

1.  $y_i = R x_i + t$
2.  $R(x_i - \bar{x}) = y_i - \bar{y} \Rightarrow y_i = R x_i - R \bar{x} + \bar{y}$

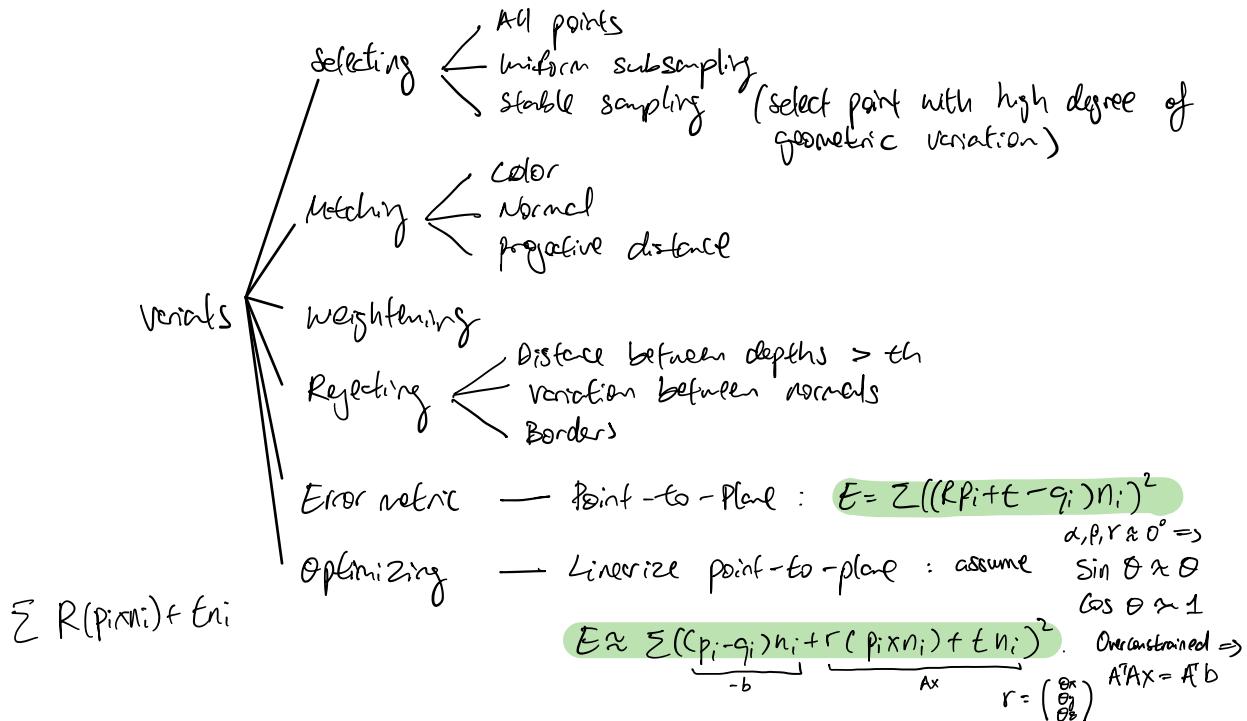
$$R x_i + t = R \bar{x} - R \bar{x} + \bar{y} \Rightarrow t = \bar{y} - R \bar{x}$$

#### ICP (no correspondence)

Iterative method that at each step approximate correspondence as closest point pairs and find alignment  $(R, t)$  based on this. Converge if consequent pairs are close enough.

Basic version:

- Select random points
- Match points on the other scan (NN)
- Reject outliers (distance  $> K \cdot \text{median}$ )
- Minimize  $E = \sum (R p_i + t - q_i)^2$



Tricks: coarse-to-fine for faster convergence and avoid local minima

## Online 3D Reconstruction

1. Capture depth map
2. Compute pose between all frames
3. Integrate depth into model (volumetric fusion)
4. Extract mesh / render SDF (raycasting)

### Volumetric fusion

weighted average of TSDF (from different views / frames)

$$D(x) = \frac{\sum w_i(x) d_i(x)}{W(x)}$$

$i = \text{camera}$

$$W(x) = \sum w_i(x) \quad \text{normalizing term}$$

$w_i(x)$  based on distance, further away lower weight

Incrementally :

$$D_{i+1}(x) = \frac{w_i(x) D_i(x) + w_{i+1}(x) d_{i+1}(x)}{w_i(x) + w_{i+1}(x)}$$

$$w_{i+1}(x) = w_i(x) + w_{i+1}(x)$$

Not true SDF:

less accurate the further away  
 more accurate with more views

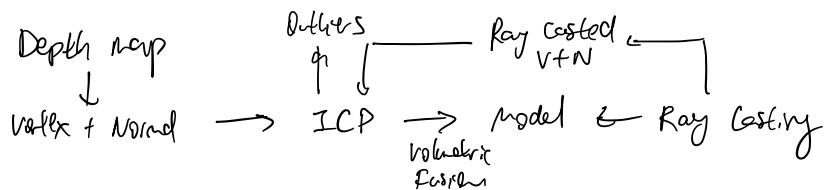
$\partial \rightarrow$  surface  
 $+$   $\rightarrow$  free space  
 $- \rightarrow$  unknown

free space covering is important

## Ray Casting

To render SDF, cast rays through focal points of each pixel.  
Traverse voxels until finding sign change (surface), compute intersection.

### Direct fusion



Issue: runs out of memory with large TSDF

### Extended Direct fusion

Rolling Volume instead of whole TSDF. Extract previous parts to mesh.

Problem: cannot revisit previous regions

### Hierarchical fusion

Voxels in tree-like structure. Efficient storage, fast lookup (ray casting)  
Problem: slow update

### Voxel Hashing

Maintain hash array (with position and pointer) and pointed data array (voxels).

Each query  $(x, y, z)$  need to be hashed (spatial hash func) to find hash entry and then with the corresponding pointer find the voxel.

Fast update, but slow look up (for ray casting)

### Point-based

Surfel representation: point + normal + radius. Problem: lower quality than SDF

## Loop closure issue

Appears when small errors (drifts) gets accumulated resulting in loop closure problem.  
Common in local optimizations like ICP of cumulative frames

## Bundle fusion

exploit surfel representation. Explicitly detect loop closure through  
localization and non-rigidly warp representation to close loop.

## Bundle fusion

Cumulative RGB-D bundle:

1. find correspondence (SIFT)

$$E(T) = w_{\text{spare}} E_{\text{spare}}(T) + w_{\text{close}} E_{\text{close}}(T)$$

$$E_{\text{spare}}(T) = \sum_{i,j} \sum_k \| T_i p_i - T_j p_j \|_2^2 \quad \text{correspondence}$$

$$E_{\text{close}}(T) = w_{\text{depth}} E_{\text{depth}}(T) + w_{\text{color}} E_{\text{color}}(T)$$

$$E_{\text{depth}} = \sum_{i,j} \sum_k \| (p_k - T_i^{-1} T_j T_d^{-1} (D_j (T_d (T_j^{-1} T_i p_k))) \cdot n_k) \|_2^2$$

$$E_{\text{color}} = \sum_{i,j} \sum_k \| \nabla I(\pi_c(p_k)) - \nabla I(\pi_c(T_j^{-1} T_i p_k)) \|_2^2$$

first optimize  $E_{\text{spare}}$  and then jointly  $E_{\text{spare}} + E_{\text{close}}$

3. Scene (model) update

Solve loop closure by de-integrating frames

$$D'(v) = \frac{w(v)D(v) - w_u(v)d_u(v)}{w(v)-w_u(v)}$$

$$w'(v) = w(v) - w(v)$$

## Non-Rigid Deformations

### Mesh Deformation

To deform meshes non-rigidly each vertex might have a different transformation.

Handlers (selected vertices) are added to the mesh, so user can drag these to desired locations.  $E_{fit}$

Regularization terms are added to move neighbor vertices (fan) accordingly.  $E_{reg}$ .

$$E(v, x) = \lambda_{fit} E_{fit} + \lambda_{reg} E_{reg}. \quad E_{reg}(v) = \sum_i \|g - v_i\|_2^2 \quad \text{over picked points}$$

### Laplacian Surface Editing

Enforce smoothness on deformations. Make neighbors move same way as 1 move.

$$E_{reg}(v, x) = \sum_i \sum_{j \in N(v)} \|(v'_i - v_i) - (v'_j - v_j)\|_2^2 \quad \text{if vertices } i, \text{ neighboring fan } j$$

### As-Rigid-as-Possible (ARAP)

Each local neighborhood has associated a different rotation (3 DoFs). Thus, each vertex has (3 local rot + 3 trans = 6 DoFs)

$$E_{reg} = E_{ARAP}(R, v) = \sum_i \sum_{j \in N(v)} \|(v_i - v_j) - R_i(v'_i - v'_j)\|_2^2 \quad \begin{array}{l} \text{Parametrize } R_i \text{ with Euler angles or} \\ \text{quaternion} \end{array}$$

$\#\text{params} = 3 \cdot n + 3 \cdot n$   
 $\#\text{constraints} = \underbrace{3 \cdot n \cdot \text{ang/fan}}_{E_{reg}} + \underbrace{3 \cdot \#\text{handler}}_{E_{fit}}$

$\left. \begin{array}{l} \text{Alternated optimization (assure one fixed)} \\ \text{vertices (translation) } \leftrightarrow \text{rotation } \Rightarrow \text{Linear (Procrustes)} \end{array} \right\}$

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

### Embedded Deformation (ED)



More DoFs than ARAP by not restricting local transformation matrix to rotations (9 vs 6 DoFs).

$$E_{reg} = E_{ED}(M, v) = \sum_i \sum_{j \in N(v)} \|(v_i - v_j) - M_i(v'_i - v'_j)\|_2^2. \quad \text{To enforce } M_i \approx R_i \text{ (orthogonality):}$$

$$E_{ref}(M) = \sum_i \text{Rot}(M_i), \quad \text{Rot}(M_i) = \underbrace{(c_1 \cdot c_2)^2 + (c_2 \cdot c_3)^2 + (c_1 \cdot c_3)^2}_{c_1 \perp c_2 \perp c_3} + \underbrace{(c_1^2 - 1) + (c_2^2 - 1) + (c_3^2 - 1)}_{c_1 \cdot c_2 = 1}$$

Harder to optimize (quadratic problem)

## Deformation Proxies

Deformation energy on proxies instead of directly on the mesh. Makes the problem sparser (less DOFs).

### Cage

(parts of) mesh tightly enclosed into a cage. Same idea as barycentric coordinates, vertex inside cage computed as linear combination of cage vertices.

$$\text{Cage } (c_1, c_2, c_3) : v_i = 0.4 c_1 + 0.1 c_2 + 0.5 c_3$$

Mesh can be deformed by animating cage.

### 3D Grid

ARAP deformation or 3D grid for realtime performance. Mesh is fitted into grids (voxels).

### Skeleton

Often used in human body reconstruction / tracking. Embed a skeleton inside the mesh, all vertices has some dependency on this joint depending on the part it belongs. Therefore, each joint has a different transformation  $M_i$ .

$$v_i = \sum_n w_{in} M_n v_i'$$

### Deformation Graphs

Applied on TD, each point takes weighted average transformation  $M_i$  from  $k$  closest graph nodes.  $n$  control nodes that affects the neighborhood.

$$p \mapsto \sum_n w_n M_n$$

## Non-Rigid Tracking

Correspondence between frames  $\Rightarrow$  movedandler.  
Find correspondence in ICP way or SIFT features.

## Non-Rigid ICP

Same as standard ICP but with more DOFs (deformation graph). Need some initial alignment

Effectively:  
find correspondence (NN / projective)  
solve for deformation and vertices

$$E = \lambda_{\text{point}} E_{\text{point}} + \lambda_{\text{plane}} E_{\text{plane}} + \lambda_{\text{color}} E_{\text{color}} + \lambda_{\text{reg}} E_{\text{reg}}$$

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

$$E_{\text{plane}} = \sum_i \|(\mathbf{v}_i - T \mathbf{v}_i') \mathbf{n}\|_2^2$$

$$E_{\text{color}} = \sum_i \|\mathcal{I}(T_i(v_i)) - \mathcal{I}(T_i(T v_i))\|_2^2$$

$$E_{\text{reg}} = \text{ARAP} = \sum_i \sum_j \|(\mathbf{v}_i - \mathbf{v}_j) - R_i(\mathbf{v}_i' - \mathbf{v}_j')\|_2^2$$

} Correspondence

Optimizing for  $v_i'$ ,  $T$ ,  
 $R_i$ .

## Template acquisition

Multi resolution template (coarse-to-fine) with prolongation operators.  
Needs a predefined template. Stored with voxel hashing

## Correspondence

Pruning: normal, distance, angle ...

Association: projective lookup, search in kernel region one block per vertex

## Optimizing

Multi resolution Gauss-Newton with internal PCG steps (linear system).

Exploit sparsity of  $J^T$  for maximal performance on GPU.

## Multi camera setup

With 8 cameras we get more constraints (correspondence), so less outliers.

## Non-Rigid Tracking and Reconstruction

Reconstruction: similar pipeline like the rigid case, Voxel Fusion.

Build initial SDF volume from depth map of 1st frame.  
In the next frames, estimate deformation from canonical undeformed model to current depth map, warp it and integrate current frame's depth map (volumetric fusion).

for each frame:

1. Take depth map
2. Extract depth map from prev. warped model  
(Marching Cubes + rasterization)
3. Align canonical model (depth map) with current frame's  $\Rightarrow$  warped model
4. Integrate current frame into canonical model
5. Next frame

### Tracking vs Reconstruction

Tracking  $\Rightarrow$  offline template reconstruction + online template tracking.  
Mitigate drift with regularizer wrt mesh.

Reconstruction  $\Rightarrow$  joint reconstruction and tracking (more dofs, more difficult).  
Template tracking wrt prev. frame  $\Rightarrow$  accumulated drift  $\Rightarrow$  reset key volume every n frames.

Challenging: RGB-D sequence under-constraint, need to calculate dofs through regularization (DL).

### Dynamic Fusion

ED (Deformation graph) + N-ICP (dense) + volumetric fusion.

Dofs determined by number of deformation graph nodes.

Limitation: in fast motion hard to determine correspondence, is the surface deforming or is new surface being scanned?

### Volume Deform

ARAP (3D grid) + N-ICP (dense depth + SIFT) + volumetric fusion

Overcome Dynamic Fusion's limitation with extra SIFT features for correspondence finding

More grid points  $\Rightarrow$  more resolution

## Real-time Geometry, Albedo and Motion Reconstruction

ED (deformation graph) + N-ICP (dense depth + dense color) + illumination correction + volumetric fusion

Optimizing the color and illumination helps with the correspondence matching in fast motion.

## Fusion 4D

More cameras (8), more complex system. Multiple canonical poses.  
Can handle some topology change.

## Motion 2Fusion

Same pipeline as Fusion 4D but improved efficiency. 200 f/s

ED (deformation graph) + key frames + multiple cameras + ML correspondences + two way ICP alignment (current  $\leftrightarrow$  reference) + Detail layer (2.5D volume)

## Lookin Good

Same pipeline. Novelty  $\Rightarrow$  neural re-rendering (DL)

## Deep Deform

Same pipeline. Novelty  $\Rightarrow$  DL learned correspondences (between current and canonical model)  
↳ heat map  
(stained hanglass)

## Body Fusion

Same pipeline but done with skeleton proxy  $\Rightarrow$  less DoFs (easier problem).

## Face Reconstruction

Done with parametric model (proxy) for dimensionality reduction purpose.  
Vertices (position) are controlled by parameters that are weights for basis functions.

$$P = (\Phi, \alpha, \beta, f, r)$$

$\Phi = [R, t]$  extrinsics of camera  
 $\alpha$  = pars corresponding to shape basis  
 $\beta$  = " altitude basis  
 $f$  = " expression basis  
 $r$  = " sh illumination basis

$$\begin{aligned} M_{geo}(d, f) &= M_{id} + E_{id} \cdot \alpha + E_{exp} \cdot f & \mu = \text{basis mean} \\ M_{alb}(\beta) &= M_{alb} + E_{alb} \cdot \beta & E = \text{basis} \\ \text{Illum}(r) &= \sum_i r_i B_i(\vec{n}) & \beta_i: \text{sh basis function}, \vec{n}: \text{point normal} \end{aligned}$$

## Illumination

Often hard (costly) problem in computer graphics  $\Rightarrow$  advanced light transport models impractical.

We use simplified models

## Environment Map

Reflection map. Panoramic image that works as a texture map for the mesh to bring illumination.

Assumptions: distant lighting, no self shadowing, no scattering

## Spherical map

## Cube map

Preferred over sphere map as it computationally cheaper, no distortion and not dependent on view point.

## Spherical Harmonics

Parametrization (approximation of env. map) with orthogonal basis defined over a sphere. With increasing number of bands (order) the approximation is more accurate. Each basis takes the normal of the point and simulates the incident lighting.

$$L = \sum_i r_i B_i(\vec{n}) \quad 3 \text{ bands} \Rightarrow 9 \text{ params (basis)}$$

This illumination model assumes a Lambertian (diffuse) surface, that is lighting is invariant to viewpoint change.

Distinct smooth lighting as env. maps.

## Parametric Model

Models constructed with a set of different face scans.

At learning time (build database/basis):

- scan faces / expressions and fit a topologically-consistent template to each with non-rigid alignment
- compute PCA-basis (set of eigenvectors) for identity, expression, albedo

At test time (fit or construct face from input image):

- find model parameters that fits input RGB(-D) face image, optimization.

## Morphable Models

Two independent PCA for identity and albedo.

More general. Parameters have global influence. Mostly caucasian people scans (bias).

## Blend shapes

Additive model controlled via blend weight. E.g. first eigenvector (basis) controls face length...

More personalized, more semantic meaning for the basis

## Principal Component Analysis (PCA)

Used for dimensionality reduction  $\Rightarrow$  compute face prior, the basis.

Given  $X = \begin{bmatrix} | & | \\ x_1 & x_2 \\ | & | \end{bmatrix}$ , where  $x_i$  are face meshes computed with non-rigid deformation over template mesh to align input scans.

1. Compute mean

$$\bar{x} = \frac{1}{N} \sum x_i \quad X_c = X - \bar{x}$$

2. Eigen decomposition (or covariance  $X_c^T X_c$  as its squared)

$$X_c^T X_c B_i = \lambda_i B_i \quad \text{with power iteration.} \Rightarrow B_i \text{ basis, } \lambda_i \text{ eigenvalue (variance).}$$

$$E_{id} = \begin{bmatrix} | & | \\ B_1 & B_2 \\ | & | \end{bmatrix}$$

$\lambda_i$  ordered decreasingly  $\Rightarrow$  leave the ones with low variance. 90% is often enough. Dim. reduction

Could also use SVD instead of PCA

## Fitting

Analysis-by-synthesis: aims to analyze a signal or image by reproducing it using a model (find parameters that synthesize the closest image possible in the span of the model).

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

$$E_{\text{sparse}}(P) = \sum \|c_i - \pi(M_i(v_j))\|_2^2 \quad \text{pixel space}$$

$$E_{\text{reg}}(P) = \sum \frac{\alpha_i}{J_{\text{id},i}} + \sum \frac{\beta_i}{J_{\text{ab},i}} + \sum \frac{\gamma_i}{J_{\text{rp},i}} \quad \text{parameters space. Helps to be } [-30, 30] \text{ to the mean.}$$

$$\left. \begin{aligned} E_{\text{geom}}(P) &= \sum \underbrace{\|M_i(P) - D(\pi(M_i(P)))\|_2^2}_{\text{point-to-point}} + \underbrace{\|[(M_i(P) - D(\pi(M_i(P)))) \cdot n_i]\|^2}_{\text{point-to-plane}} \\ E_{\text{color}}(P) &= \sum \|M_i(P)_c - I(\pi(M_i(P)))\|_2^2 \end{aligned} \right\} \begin{array}{l} \text{over rendered} \\ \text{face mask.} \\ \text{But in world} \\ \text{space} \end{array}$$

## Differential rendering

Dense terms residuals are calculated in pixel space, but we want the partial derivatives to affect the vertices/colors in world space.

Each pixel in the rendered image comes from a triangle in world space with location defined by the barycentric coordinate inside the triangle. If we use (store) this information, then we are able to calculate the derivatives.

$$C_{\text{pixel}} = \alpha \cdot C_{P_i} + \beta \cdot C_{P_j} + \gamma \cdot C_{P_k} \quad (\alpha, \beta, \gamma) \text{ barycentric coordinate.}$$

The color is just a linear combination of the triangle's vertices' colors.

## Tricks

The computation cost comes from the dense color terms.

To make it work in real-time we can adopt a coarse-to-fine strategy which not only helps the speed but also the convergence (smoother function). The result from a lower resolution optimization are used as initialization for next level.

## RGB Only

Also doable with RGB only by dropping the Ego terms. But need several different frames to jointly optimize for the identity and albedo (base color, shape, texture). Then, these are fixed and only expression, extrinsics and illumination are optimized.

The mesh would not look good, but the rendered image is still ok.

## Mouse, Teeth

Use proxies: add generic teeth mesh inside mouth. Depends on RGB-D

An improved solution (RGB only) is to save the square cropped mouth + teeth and compute a descriptor for each. Then after the expression transfer, use parameters such as expression ( $\delta$ ), landmarks, R and define some similarity metric to find the best matching frame.

## Body and Hand

Same approach as faces, parametric model.

- Take a set of scans from variety of subjects / poses
- Non-rigidly align them with template for topology consistency.

## Body

pose defined from embedded skeleton.

$$n(\vec{\theta}, \vec{\beta}, \vec{u}, \vec{\Phi})$$

$\vec{\theta}$  = pose (skeleton)  
 $\vec{\beta}$  = shape  
 $\vec{u}$  = texture  
 $\vec{\Phi}$  = hyper params

## CAESAR

One of the first human body scan database.

125 male + 125 female all neutral pose.  
74 markers (landmarks) to help the alignment.

Can't handle pose

## SCAPE

Expansion of CAESAR with more identity and poses.  
37 identity + CAESA, 70 poses of a particular person.

Parametric (PCA) model for base shape + Rigid (skeleton) and Non-rigid muscle flex.  
 $P = (R, \beta, Y)$  alternating optimization.  $Q$  (corresponding to muscle flex) can be obtained using joint shape predicted mesh the given  $Q$  and optimized  $R$ .

## SMPL

Novel way to parametrize pose by skeleton. Separate human body by parts controlled by skeleton joint's movement [Bt].

The rotation can be represented with 3x3 mat then adding constraints to enforce orthogonality ( $\approx$  ED) or better with exponential maps ( $SO(3)$ ) to reduce dimensionality.

$\|\vec{w}_j\|$  angle of rotation  $\vec{w}_j$  axis of rotation.

$$R = e^{\hat{\vec{\omega}}} = I + \hat{\vec{\omega}} \sin(\|\vec{w}_j\|) + \hat{\vec{w}}^2_j (1 - \cos(\|\vec{w}_j\|))$$

Get back to lie group  $SO(3)$  from skew-symmetric  $\hat{\vec{\omega}} \in so(3)$

So each joint's pose now is represented by

$$G(\vec{\omega}; j) = \begin{bmatrix} \vec{c}_{3 \times 3} \\ 0 \end{bmatrix}$$

The total motion of a vertex (or part) is affected by the concatenation of all joint's movement:

$$\vec{p}_s = G(\vec{\omega}_1, \vec{\omega}_2, j_1, j_2) = G(\vec{\omega}_1, j_1) G(\vec{\omega}_2, j_2) \vec{p}_b \Rightarrow$$

$$\vec{t}_i' = \sum_k w_{ki} G^k(\vec{\theta}, J) \vec{t}_i \quad k = \text{joints for smooth transition (blended linear comb.)}$$

$w_{ki}$  = weight that a vertex has for a joint  $k$

$$M(\vec{\beta}, \vec{\theta}; \Phi) : \mathbb{R}^{10N \times |\beta|} \rightarrow \mathbb{R}^{3N} \quad \text{the shape model} . \quad \vec{\beta} : \text{shape params, } \vec{\theta} : \text{joint rotations}$$

$W(\vec{T}, J, W, \vec{\Theta}) : \mathbb{R}^{?} \rightarrow \mathbb{R}^{3N}$  gives the displacements of each vertex.

$\vec{T}$  = template mesh,  $J$  = joints,  $W$  = weights,  $\vec{\Theta}$  = joint rotation transf.

Training  $\Rightarrow$  find  $\Phi = (\vec{T}, W, J, S, P)$ . Inference  $\Rightarrow$  find  $\vec{\beta}, \vec{\theta}$   
the model.  $S$  shape basis  
 $P$  (joint) pose basis

Artifacts (unnatural local deformation) from linear blend skinning, as it's just a linear function.  $\Rightarrow$  Muscles should flexed as pose changes.

Solution: a blend shape with a set of vertex displacement relative from the neutral pose (offsets)

When the pose is changed, the base shape should also be changed in a small degree (offset)

## Reconstruction

Fit template to input scans (SMPL model):

$$E = W_d E_d + W_s E_s + W_m E_m$$

$$E_d = \sum_i w_i \cdot \text{dist}^2(T_i, V_i, D) \quad \text{data term (2D silhouette / 3D)}$$

$$E_s = \sum_{(i,j) \in \text{Edge}} \|T_i - T_j\|_F^2 \quad \text{smoothness regularization term}$$

$$E_m = \sum_i \|T_{Vi} V_{Ui} - M_i\|^2 \quad \text{marker term (2D)}$$

Possible with RGB only with silhouettes and segmentation...

## Hand

Similar idea.  $\beta$  = shape  $\theta$  = pose. However dataset are usually built with synthetic data as scanning hands is usually difficult.

- 2 approaches
- Discriminative: per-frame reinitialization to prevent error propagation
  - Generative: template mesh deforming to align with captured depth data

### Taylor et al. 16

Subdivision model which is mostly  $C_2$  continuous (twice differentiable)  $\Rightarrow$  fast LM fitting.

Fit Kinect input data to sparse points of hand model.

Blend shape + linear blend skinning for pose.

At each step:

1. Preprocess
2. Generate starting point from last frame's prediction (reinitializer)
3. Optimize all  $E(\theta)$
4. Take the one that yields lowest cost

$$E(\theta) = E_{\text{data}} + \lambda_{\text{bg}} E_{\text{bg}} + \lambda_{\text{pose}} E_{\text{pose}} + \lambda_{\text{light}} E_{\text{light}} + \lambda_{\text{tmp}} E_{\text{tmp}} + \lambda_{\text{inf}} E_{\text{inf}} + \lambda_{\text{tips}} E_{\text{tips}}$$

### Sphere-mesh (Tkach)

sphere proxies for the hand  $\Rightarrow$  articulated and flexible components.

Hybrid approach: discriminative (estimate the template) + generative (refine results, real-time tracking)

Similar energy formulation than Taylor et. al.

Posterior work: added KFBA filter

### Issue of parametric model

Need expensive scanning setup to build database (PCA)

Neutral pose is not identical between subjects

Ambiguity to distinguish identity or pose change (need regularizers)

Drift in non-rigid registration (useless PCs)