

Render - Register - Deform

1) Deep CPD [Category level 3D non rigid ...]

→ Imp concept [deformation fields] for visible parts
+ shape [latent] space for inducing deformation of occluded part

→ flow net 2 [multiplied with 2RGB, 2 mask]

\downarrow out
[3 channel for xyz for each point]

→ coherent point drift [CPD]

2 sets D-dimensional X, Y

CPD $\xrightarrow{\text{output}}$ deformation field mapping Y towards X

so we consider Y points as centroids of Gaussian Mixture Model

X points are considered as samples drawn from these GMMs

- equal isotropic covariances σ^2 and equal membership probabilities

$$P(m) = \frac{1}{m} \text{ for all GMM component}$$

- CPD outputs \uparrow deformed a deformed pointset as maximization of the probability of X being drawn from GMM, by moving centroids Y in a coherent manner [Expectation Maximisation problem]

$$\therefore \pi \Rightarrow \pi(Y, v) = Y + v(Y)$$

\downarrow
initial pos.

\nearrow Displacement function

coherent point drift for any D-dimensional set of points $Z \in \mathbb{R}^{N \times D}$ v is defined as

$v(Z) = G(Y, Z)W$, $W \in \mathbb{R}^{N \times D}$

$G(Y, Z)_{ij} = \exp\left(-\frac{1}{2\beta^2} \|y_i - z_j\|^2\right)$

\downarrow
estimated in
M step of EM algo

\nearrow D-dimensional deformation vector

- method

$C \in \mathbb{R}^{n \times 3}$ point cloud

$C_m \in \mathbb{R}^{m \times 3}$ mesh

$$C'_m = C_m + G(C_m, C) \underbrace{W(C, D)}_{\substack{\hookrightarrow \text{offsets} \\ \text{to mesh}}}, W(C, D) \in \mathbb{R}^{n \times 3}$$

\nearrow offset of y instance
 $D = \text{observed offsets to } C$
 \nearrow deform towards

$$w(c, 0) = g(c, c) w(c, 0)$$

offset matrix represented as
a 3 channel image $\times_{1,0,2}$

- deformation on occluded parts

[shape space] for a known object category

$$\tilde{r}_i(c, w_i) = c + g(c, c) w_i$$

\hookrightarrow lower dimensional space spanned by principle components of all deformation fields w_i , $w_i \in \mathbb{R}^{3n}$

\therefore we can apply PCA Expectation Maximization [PCA-EM]
to find latent space of category.

\therefore mat $L \in \mathbb{R}^{L \times 3n}$ containing L 's principle components is determined and point $x \in \mathbb{R}^L$ in latent space can be mapped into feature vector $\tilde{w} \in \mathbb{R}^{3n}$ by $\tilde{w} = Lx + \bar{w} \rightarrow$ mean of all feature

\therefore Canonical Mat C with L represent 3D latent model or object latent vectors

So find \tilde{w} for all points of C occluded

$$\therefore \hat{s} = \bar{g}(c, c)(Lx + \bar{w})$$

\downarrow
 $\mathbb{R}^{3n \times 3n}$ same as $g(c, c)$ but with added zeros to match space

Main loss is \therefore

$$L(\delta_{vis}, \hat{s}) = \sum_{v \in vis} \| \delta_{vis}(v, \cdot) - \hat{s}(v, \cdot) \|^2$$

\downarrow
visible indices

for training inverse deformation is used based on interpolation

2) Iterative 3D Seg

\rightarrow optimization in Image loss

- using U-net features $I_{out} - I_{out}$

- ~~6D~~ 6D Pose estimation

3) Deformable Seg [Kov et al]

$$\Pi: \mathcal{C}^3 \rightarrow \mathcal{C}^2 \quad \text{pin hole camera}$$

$$M_0 \rightarrow CT$$

M^* → seg true volume

$$\text{Register } (M_0) \rightarrow M \quad \downarrow \text{evolved model}$$

using isotropic Neo-Hookean elastic Model

ϕ_{det} per tetrahedron

$$\text{contour constraint } \phi_{\text{contour}}(m) = \text{dist}(c, c^*) \quad C \in \mathbb{R}^2$$

∴ Silhouette $C\Pi(m)$ etc

Ridge $C\Pi(m)$ gaussian curvature operator K auto or selectable by user
utilizing all

$$\phi_{\text{point}}(m) = \text{dist}(n, L^*) \quad n \in m \text{ end point of contour seg } C \\ L^* \text{ sight line passing through this points} \\ \text{image pos on } C^* \text{ ridge-segment}$$

- Uses shading to better adapt the model M' to Image I'

- visible surface emerges as a triangle mesh in the tessellated M'

∴ ϕ_{shading} for each triangle [compare grey level of image pix with projected triangle
 $\phi_{\text{shading}}(m) = p(n, m^*)$ ~~triangle~~ projected triangle

- optimisation

- 1) contour dist optim hist
- 2) new shading optim

$$\begin{aligned} \phi: S &\rightarrow \mathcal{R} \\ \mathcal{R} \text{ solution set} & \quad \Pi: S \rightarrow \mathcal{L} \end{aligned}$$

$$\cancel{\text{def}} \quad M \in \arg \min_M \sum_{i=1}^m \text{dist}(M, s_i)$$

Basically search M' to all constraint solution sets
See fig 1 for optim pipeline

-) constraint mappings are linearized version of ϕ_{det}

4) Learning deformable Tet Mesh for 3D scenes [DefTet]

~~Color and visibility gradient back projected~~

-) 3D spatial graph that tetrahedralised
-) Each tet is a polyhedron with 4 triangular faces, ~~6~~ edges and 4 vertices. Neighboring tet share ~~one~~ face, 3 edges, 3 vertices
-) So how to form a 3D object:
 - Predict offsets for each vertices of the initial tet
 - Predict occupancy of each tet

$$\cdot) \mathbf{v}_i = [x_i, y_i, z_i]^T \text{ vertex of tet}$$

$$f_s = [v_{a_s}, v_{b_s}, v_{c_s}] \text{ Triangular face of tet}$$

$$T_k = [v_{a_k}, v_{b_k}; v_{c_k}, v_{d_k}] \text{ Tet}$$

-) Binary occupancy of tet T_k as O_k

$$O_k = O_k(e(T_k)) \quad \text{func } e \text{ derives different occupancy functions
e.g. occupancy wrt tet's centroid}$$

-) $P_s(v) = O_{s_1}(v_{s_1})(1 - O_{s_2}(v_{s_2})) + (1 - O_{s_1}(v_{s_1}))O_{s_2}(v_{s_2})$
 - probability that a triangular face f_s derives a surface of an object
 - T_{s_1}, T_{s_2} neighboring tets that share face f_s
 - prob that neighboring tet to be occupied and the other not

$$\left[[\Delta x_i, \Delta y_i, \Delta z_i]_{i=1}^N, [O_k]_{k=1}^K \right] = h(v, I; \theta)$$

I = image or point cloud

h = neural net → to predict ~~deformation~~ and occupancy of tet

θ = parameters N·N h offset of vertex

∴ deformed verti are

$$V'_i = [x_i + \Delta x_i, y_i + \Delta y_i, z_i + \Delta z_i]^T$$

-) Inference ~~of~~ of the N·N is deformed tets and occupancy when given I .

) Total loss :

$$L_{\text{all}}(\theta) = \lambda_{\text{recon}} L_{\text{recon}} + \lambda_{\text{vol}} L_{\text{vol}} + \lambda_{\text{lap}} L_{\text{lap}} + \lambda_{\text{sm}} L_{\text{sm}} + \lambda_{\text{delta}} L_{\text{delta}} \\ + \lambda_{\text{amps}} L_{\text{amps}}$$

L_{recon} = reconstruction loss that aligns surface faces with ground truth faces depending on 2D or 3D.

L_{lap} = Laplacian loss to penalize change in relative positions of neighboring verts

L_{delta} = penalizes large deformations

L_{vol} = equiv volume loss penalizes difference in volume for each tet

L_{unif} = penalizes distortion of each tet

L_{sm} = smoother loss penalizes difference in normals for neighboring faces when 3D G.T available

) Recon loss with 3D G.T

$$L_{\text{occ}} = \sum_{k=1}^K \hat{\Omega}_{ik} \log(\Omega_{ik}) + (1 - \hat{\Omega}_{ik}) \log(1 - \Omega_{ik}),$$

- Supervise occupancy prediction from 'h' with Binary Cross Entropy

- Ω_{ik} ground truth occupancy for Tet T_k

- If winding number of centroid of a tet is positive consider it occupied

) Deformation :

obtain G.T occupying whose $P_S(v)$ equals 1, occupied faces as F. Deform F to align with target objects surface.

$$L_{\text{surf}}(\theta) = \sum_{p \in S} \min_{f \in F} \text{dist}_f(p, f) + \sum_{q \in S_F} \min_{p \in S} \text{dist}_p(p, q)$$

$\text{dist}_f(p, f)$ = L2 dist from point p to f

$\text{dist}_p(p, q)$ = L2 dist point p to q

$$L_{\text{recon3D}}(\theta) = L_{\text{occ}} + \lambda_{\text{surf}} L_{\text{surf}}$$

i) Recon via depth rendering

-) Recon when 2D available
-) Algo to project deformable tet to a 2D image using the given camera.
-) Consider initial Mesh with V' and F' .
- i) each vertex v_i with an RCB attribute C_i and visibility D_i ;
- i) shoot ray from ~~pixel~~ camera origin to pixel I_j to a face F' in the mesh. $[f_{L_1}, \dots, f_{L_L}]$ L total number of rays and the faces are ordered in increasing depth.
- i) Render all faces and softly combine via visibility.
- i) $F_{lk} = (v_a, v_b, v_c)$, attribute in pixel I_j is obtained via barycentric face

$$[w_a, w_b, w_c] = \text{Barycentric}(v_a, v_b, v_c, I_j)$$

$$D_j^k = w_a D_a + w_b D_b + w_c D_c, C_j^k = w_a C_a + w_b C_b + w_c C_c$$
- i) $m_k = \prod_{i=1}^{k-1} (1 - D_j^i) D_j^k$ visibility m_k of face F_{lk} in pixel I_j depends on its own visibility and all the faces in front of it.
- i) $M_j = \sum_{k=1}^L m_k, R_j = \sum_{k=1}^L m_k C_j^k$
- * for pixel I_j we use a weighted sum to obtain color R_j and visibility M_j
- * Based on chain rule gradients from M_j and R_j can be naturally backpropagated to vertex positions
- * They found that per vertex occupancy and choosing the Max works well
 $\therefore O_k = \max(D_{a,k}, D_{b,k}, D_{c,k}, D_{d,k})$

i) 2D Loss:

$$L_{\text{depth 2D}} = \sum_{j=1}^J |R_j^{gt} - R_j|_1 + \lambda_{\text{mask}} |M_j^{gt} - M_j|_1,$$

M^{gt} is the mask for object in Image

- 1) Network architecture from DVR [differentiable volumetric rendering...]
 training with 3D or 2D G.T
 - Here the output of DVR the last layer is changed to accommodate tetrahedron shape. Initial tet of $4 \times 4 \times 4$

t_θ = texture/color Network

f_θ = occupancy Network

$$L(\hat{I}, I) = \sum_u \|\hat{I}_u - I_u\| \quad I_u = \text{RGB val of ray}$$

$$\frac{\partial L}{\partial \theta} = \sum_u \frac{\partial L}{\partial \hat{I}_u} \cdot \frac{\partial \hat{I}_u}{\partial \theta}$$

since $\hat{I}_u = t_\theta(\hat{p})$ ray intersection point (\hat{p})

$$\therefore \frac{\partial \hat{I}_u}{\partial \theta} = \frac{\partial t_\theta(\hat{p})}{\partial \theta} + \frac{\partial t_\theta(\hat{p})}{\partial \hat{p}} \cdot \frac{\partial \hat{p}}{\partial \theta}$$

$$\frac{\partial \hat{p}}{\partial \theta} = \frac{\partial r(\lambda)}{\partial \theta} = w \frac{\partial \lambda}{\partial \theta} \quad r(\lambda) = r_0 + \lambda w \vec{v}_{\text{rec}}$$

λ = surface depth

$$\frac{\partial \hat{p}}{\partial \theta} = - \left(\frac{\partial f_\theta(\hat{p})}{\partial \hat{p}} \cdot w \right)^{-1} \frac{\partial f_\theta(\hat{p})}{\partial \theta}$$

5) Fast gradient descent for surface capture via diff rendering

- .) Non stochastic CR-D use
- .) Analogous to plexroducts
- .) Extract 3D mesh
- .) Volumetric $\underbrace{R \cap B A}_{\text{color opacity}}$ like res

.) Input : In \rightarrow multiview Calib images

$B_n \rightarrow$ Background images [captured before object is placed in scene]

- .) Use diff volume recon with appropriate regularization terms
- .) opacity parameterisation could be voxel, ~~SH~~, FCNN
- .) RIBA, S-H, NN etc.

All optimization of volumetric diff rendering is based on color similarity loss

$$L_{\text{photo}} = \sum_{r \in \text{rgs}} \| C(r) - C_{\text{gt}}(r) \|$$

\downarrow sampled \downarrow gt color

~~defn~~ $\alpha_i \in [0, 1]$ transparency
 $c_i \in [0, 1]^3$ color

.) Sparse grid representation where RIBA obtained by interpolation or in plixels.

$$T_K = \prod_{j < K} \alpha_j \quad \text{partial transmittance} \quad K \text{ is a value}$$

$$C_K = \sum_{j \leq K} (1 - \alpha_j) c_j T_j \quad \text{partial color} \quad \text{regular intervals}$$

$$C = \sum_{K < J} (1 - \alpha_K) C_K T_K + T_\infty C_\infty$$

.) color and transparency of a sample obtained by linear interpolate of 8 closest voxels $\alpha = \text{lerp}(\bar{\alpha})$ $\bar{\alpha}$ bar indicates voxel val

$$\text{view dependent effects obtained by S-H} \quad \begin{matrix} & \nearrow \text{bands} \\ [\text{distance dependent highlight}] & \end{matrix}$$

$$\bar{C}(v) = \sum_{l=0}^B \sum_{m=-1}^L Y_{l,m}(v) \beta_{l,m}$$

$Y_{l,m}(x_1, y_1, z)$ func
are polynomials

color control: $\beta_{l,m} \in [-1, 1]^3$

$v =$ unit dir
of ray

.) losses

1) Robust photometric loss

$$\|n\|_H = n^2 \text{ if } n < \epsilon \text{ else } \epsilon(|n| - \frac{\epsilon}{2}) \quad \text{Huber norm}$$

$$\therefore L_{\text{photo}} = \sum_{\text{rays}} \|C(r) - C_{\text{gt}}(r)\|_H$$

.) Spatial Regu

$$L_{\text{smooth}} = \sum_{\text{voxels}} \|\nabla \alpha\|_2^2 \quad \begin{array}{l} \text{Total variation} \\ \text{regularization} \end{array}$$

.) Ray color consistency

$$L_{\text{sample}} = \sum_{\text{sample}} \|c_k - C\|_H T_k(1-a_k) \quad \begin{array}{l} \text{favor voxel color} \\ \text{consistency} \end{array}$$

.) S.H

$$L_{\text{S.H.}} = \sum_{\text{voxels}} \|\beta_{l,m}\|_H$$

.) Ballotary occupancy bias

$$L_{\text{ballot}} = \sum_{\text{voxels}} \alpha^{-2}$$

for difficult regions with no texture or hard to reconstruct, this promotes more opaque voxels

.) Color normalisation

$$C_{\text{gt}} = g C_{\text{uncorrected gt}} + b$$

$$C_{\text{bg}} = g C_{\text{uncorrected bg}} + b$$

]} Camera gain and bias for different exposure

.) Gradi descent

we have analytical grads here

$$\frac{\partial L_{\text{photo}}}{\partial a_k} = \frac{\partial L_{\text{photo}}}{\partial C} \frac{\partial C}{\partial a_k}$$

$$\frac{\partial L_{\text{photo}}}{\partial C} = \|C - C_{\text{gt}}\|_H'$$

$$\frac{\partial C}{\partial a_k} = -C_k T_k + \frac{C - C_k}{a_k}$$

) March 1st ray to build T₀ and C then second ray to build $\frac{\partial L_{photo}}{\partial \alpha_k}$

$$\frac{\partial L_{photo}}{\partial \bar{\alpha}} = \sum_{r \in rays} \sum_{k \in S(r)} \frac{\partial L_{photo}}{\partial \alpha_k} \frac{\partial \alpha_k}{\partial \bar{\alpha}}$$

) We don't take all vals in $S(r)$:: voxel of interest contributes
 $S(r)$ = set of samples on ray 'r' to which we interpolate per sample gradients

$$\frac{\partial L_{photo}}{\partial \bar{\alpha}} \approx \sum_{\text{camera}} \text{lens} \left(\frac{\partial L_{photo}}{\partial \alpha} \right)$$

) Surface extraction done using transparency vals $\bar{\alpha}$ using marching cubes, color is computed per pixel at rasterisation time.

~~6) Using Differentiable Tetrahedral Meshes for 3D reconstruction~~

6) Deep Marching Tetrahedron: A hybrid representation for high Res 3D shape synthesis [DMTET]

- ~~1) Optimizes for voxel placement and occupancy
is differentiable wrt. 3D recon losses~~
- 1) DMTet directly optimizes for the reconstructed surface, which allows better detail formation.
- 2) Contributions:
- 1) Using Marching Tet as a differentiable iso-surface layer allows topological change for the shape shown by an implicit field.
 - 2) Incorporate MT in a DL framework, combining implicit + explicit fields.
 - 3) Coarse to fine optimization to scale DMTet to higher Res
- 3) Our Network learns grid deformations and subdivision jointly to better represent surface! Not relying on pre-computed hierarchy.
- 4) Deep implicit field [DIF] represent a 3D shape as a zero level set of a continuous function parameterized by a NN.
- 5) Optimized voxelised 3D input to a mesh that incorporates G.T. data [course voxel]

1) Extracting triangular 3D Models [Nvdiffrc]

$$\underset{\phi}{\operatorname{argmin}} \mathbb{E}_c [L(I_{\phi}(c), I_{\text{ref}}(c))]$$

Empirical Risk

ϕ = para [SDF, view offset, etc]

$I_{\phi}(c)$ = Image from renderer
'c' pose

$I_{\text{ref}}(c)$ = ref image

L = loss

.) Transparency/translucency not supported

.) Usually works with 50+ views assuming 360 coverage
requiring pose etc just like ref

.) Main contrib is the mipmap differentiable Lighting