

Neural Human Rendering

Recent novel view synthesis method of dynamic human performance



Sogang University

Vision & Display Systems Lab, Dept. of Electronic Engineering



Presented By

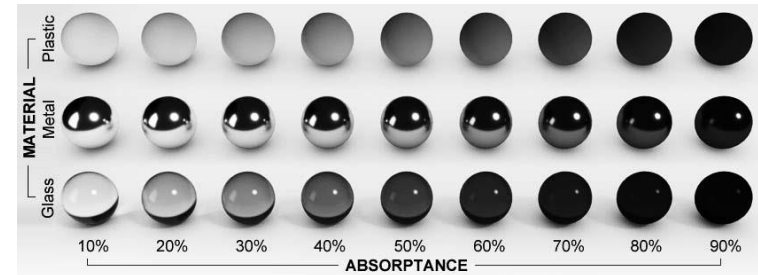
Hosung Son

Contents

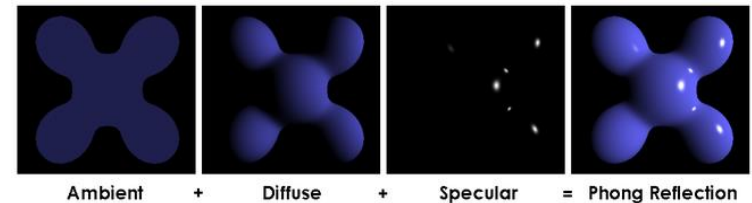
- What is neural rendering?
 - Computer graphics in 3d rendering
 - Neural rendering
- Fundamentals of Neural Rendering
 - Key points of Neural rendering
 - Scene representation
- Novel view synthesis
 - Static contents
 - Non-static contents
- Paper
 - HumanNeRF: Free-viewpoint Rendering of Moving People from Monocular video

What is neural rendering?

- Computer graphics in 3d rendering
 - Physical parameters from object and camera
 - Light transport
 - ☼ Absorption
 - ☼ Reflection
 - ☼ Scattering
 - Material properties
 - Camera parameters for image projection
 - All parameters should be input rendering model for high-quality reconstruction
 - Rendering equation¹⁾
 - Consider only emitted, scattered light and geometry
 - More considerations are in next version of the eq.
 - ☼ Reflection, Transmission...
 - Ray tracing, ray marching, path tracing



Visual differences according to light absorptance



Phong shading model

$$I(x, x') = g(x, x') \left[\epsilon(x, x') + \int_S \rho(x, x', x'') I(x', x'') dx'' \right]$$

where:

- $I(x, x')$ is related to the intensity of light passing from point x' to point x
- $g(x, x')$ is a "geometry" term
- $\epsilon(x, x')$ is related to the intensity of emitted light from x' to x
- $\rho(x, x', x'')$ is related to the intensity of light scattered from x'' to x by a patch of surface at x'

Rendering equation

What is neural rendering?

- Computer graphics in 3d rendering

- Surface rendering

- MVS with SfM [COLMAP]

- ⚡ Feature extracting/matching algorithm
 - ⚡ Triangulation

- MVS with Neural Networks

- ⚡ Feature extraction using DNN
 - ⚡ Matching cost volume (Homography)
 - ⚡ Depth estimation per view (Regression)

- Volume rendering

- Based on ray casting method

- 1) Casting rays
- 2) Sampling
- 3) Shading
- 4) Compositing

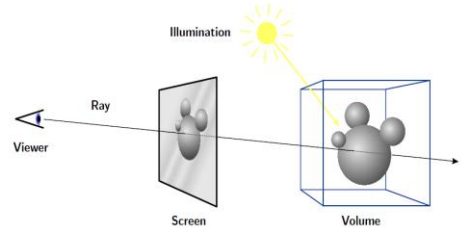
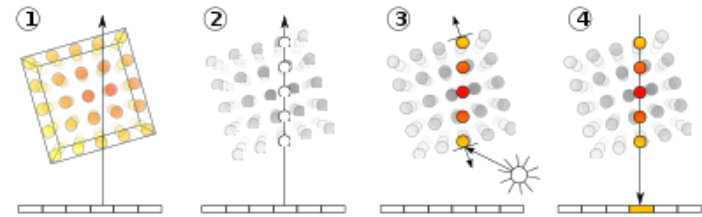
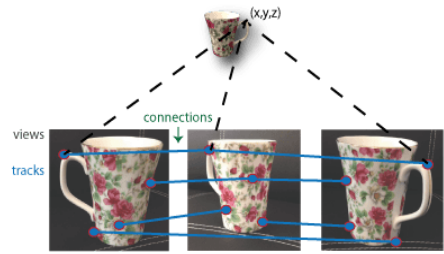


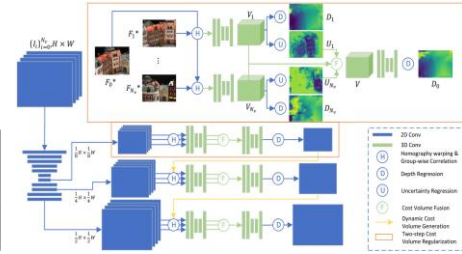
Illustration of ray casting



(1) Ray casting, (2) Sampling, (3) Shading, (4) Compositing



Structure-from-Motion



MVS Network architecture



Family



Panther

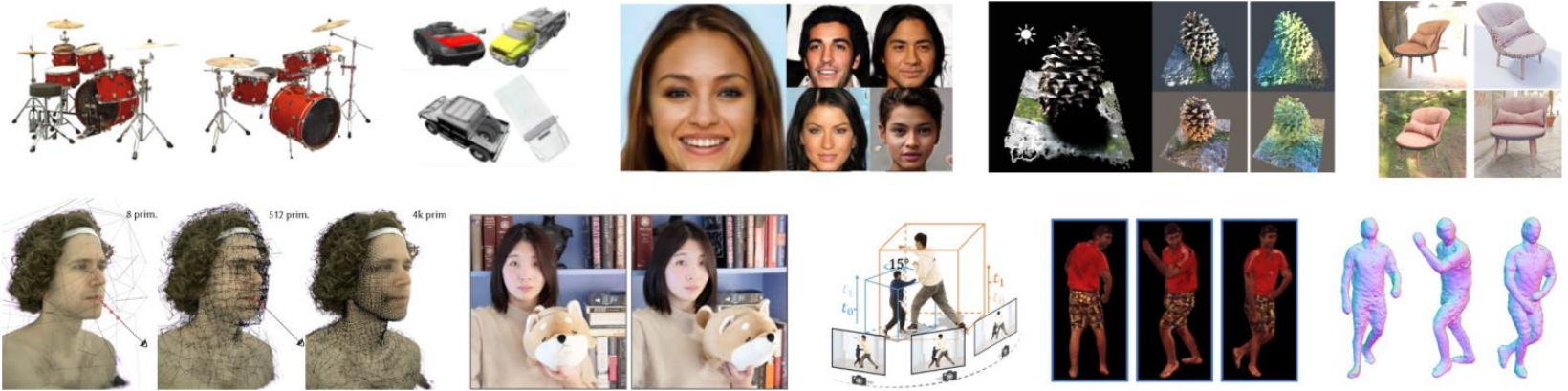


Horse

Point cloud results of Vis-MVSNet

What is neural rendering?

- Neural rendering



Many kinds of neural rendering tasks

- Shape and appearance rendering combining two insights.
 - Classical Computer Graphics
 - Deep Neural Network
- Optimizing functions with Neural network such as MLPs
 - Non-linear optimization
 - Optimization strategies

Fundamentals of Neural Rendering

- Key points of Neural rendering
 - Disentanglement of camera capturing process and 3D scene representation
 - Neural rendering methods should be differentiable for training

- Scene representation

- Surface representation

- Explicit → point cloud, polygon mesh

$$\ast S_{explicit} = \left\{ \left(\begin{array}{c} x \\ y \end{array} \right) \mid \left(\begin{array}{c} x \\ y \end{array} \right) \in \mathbb{R}^2 \right\}$$

- Implicit → zero level set of implicit function

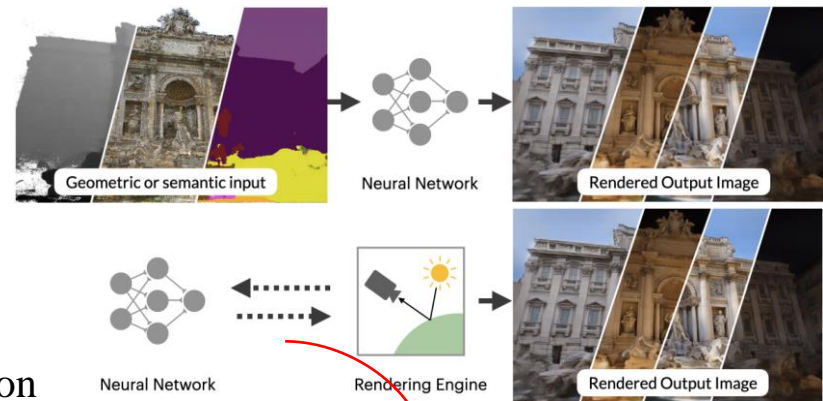
$$\ast S_{implicit} = \left\{ \left(\begin{array}{c} x \\ y \\ z \end{array} \right) \in \mathbb{R}^3 \mid f_{implicit}(x, y, z) = 0 \right\}$$

- Volumetric representation

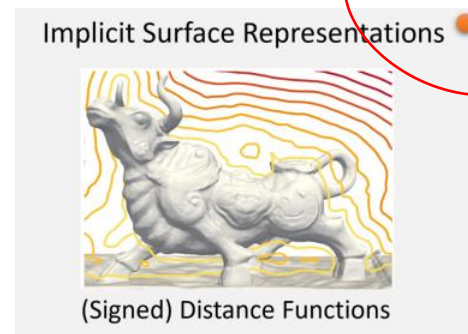
- Densities, opacities or occupancies

- Multi-dimensional features

- ∴ colors, radiance



Difference between 2D and 3D neural rendering

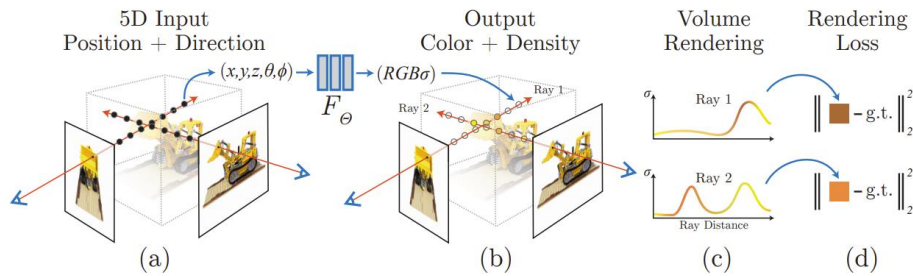


differentiable scheme from computer graphics which are motivated by physics.

Implicit surface representation of SDF

Novel view synthesis

- View synthesis of static contents
 - Neural Radiance Fields (NeRF¹⁾)



- Volume rendering with radiance fields²⁾

$$C(\mathbf{r}) = \int_{t_n}^{t_f} T(t) \sigma(\mathbf{r}(t)) \mathbf{c}(\mathbf{r}(t), \mathbf{d}) dt, \text{ where } T(t) = \exp\left(-\int_{t_n}^t \sigma(\mathbf{r}(s)) ds\right)$$

- Uniform ray sampling

$$t_i \sim \mathcal{U}\left[t_n + \frac{i-1}{N}(t_f - t_n), t_n + \frac{i}{N}(t_f - t_n)\right]$$

- Discretized representation of volume rendering³⁾

$$\hat{C}(\mathbf{r}) = \sum_{i=1}^N T_i (1 - \exp(-\sigma_i \delta_i)) \mathbf{c}_i, \text{ where } T_i = \exp\left(-\sum_{j=1}^{i-1} \sigma_j \delta_j\right)$$

Volume rendering eq. from NeRF¹⁾

$C(\mathbf{r})$: expected color
 $\sigma(\mathbf{x})$: volume density
 $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$: camera ray
 t_n : near bound
 t_f : far bound

$\delta_i = t_{i+1} - t_i$: distance between samples

$$C(\mathbf{r}) = \sum_{i=1}^D \left(\prod_{j=1}^{i-1} (1 - \alpha_j)\right) \alpha_i \mathbf{c}(\mathbf{x}_i),$$

$$\alpha_i = 1 - \exp(-\sigma(\mathbf{x}_i) \Delta t_i),$$

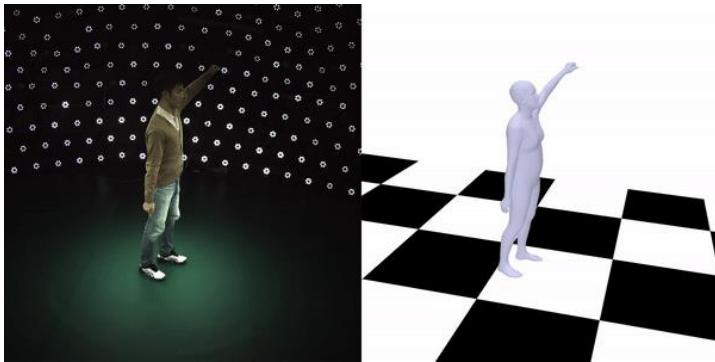
Volume rendering eq. from HumanNeRF

Quadrature rule

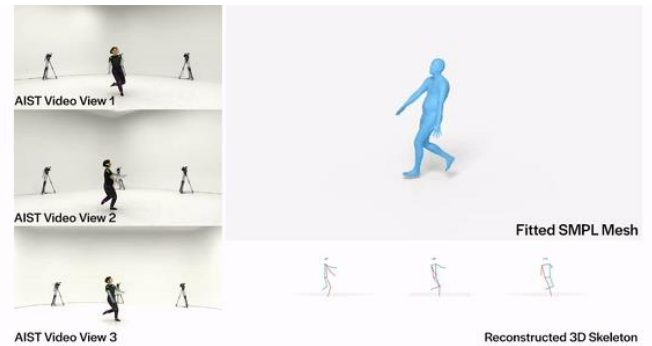
→ Differentiable!!

Novel view synthesis

- View synthesis of Non-static contents
 - Human performance dataset (Multi-view)



ZjU-MoCaP¹⁾



AIST++²⁾

- “In the wild” monocular videos from Youtube



story



invisible



way2sexy

Novel view synthesis

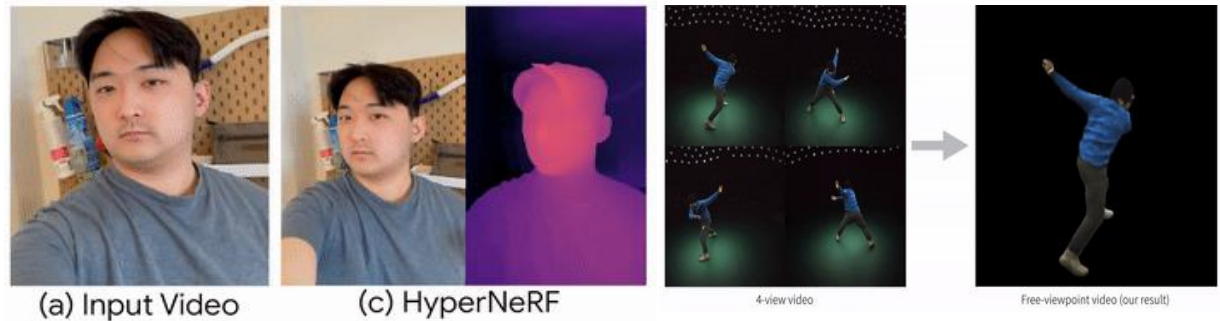
- View synthesis of Non-static contents

- Time varying in motion

- Since human performance is time varying, it is difficult to generate generalized appearance.
- If the motion is highly dynamic, unintentional artifacts could be input to model.

- Deformation

- Human face
 - ⊛ expression
- Clothes
- Etc...



(a) Input Video

(c) HyperNeRF

4-view video

Free-viewpoint video (our result)

NeRF with face deformation [HyperNeRF¹]

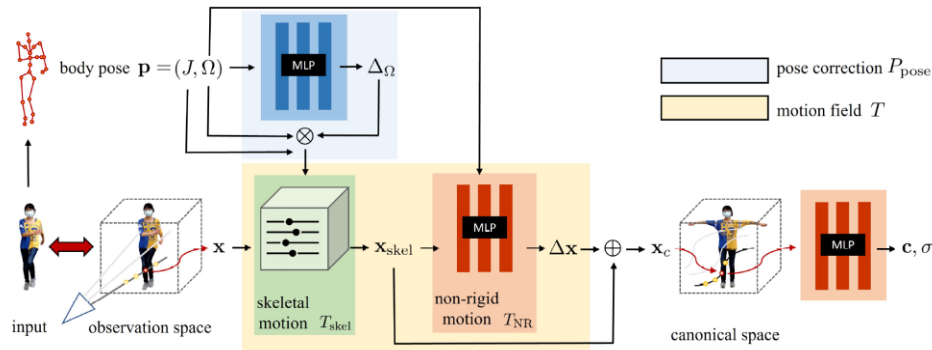
NeRF with dynamic motion [NeuralBody²]

- Occlusion

- If not a multi-view contents, occlusion is very likely to occur.
- Cross section of human body parts

Paper – HumanNeRF [CVPR2022]

- Overall architecture



HumanNeRF teaser video

- Input

- Monocular video of complex human performance
- Human, camera pose at each frame (Not use template)

- Output

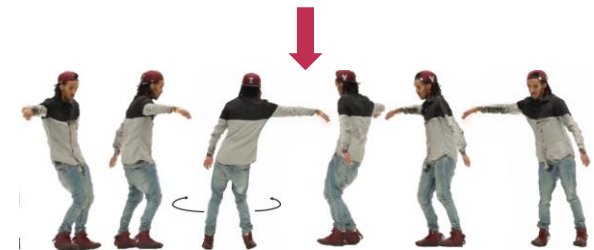
- Free-viewpoint rendering for any frame in the sequence

- Components

- Pose refiner: $\Delta\Omega$
- Motion fields: $T_{\text{skel}}, T_{\text{NR}}$
- Canonical Volume: F_c



Monocular video



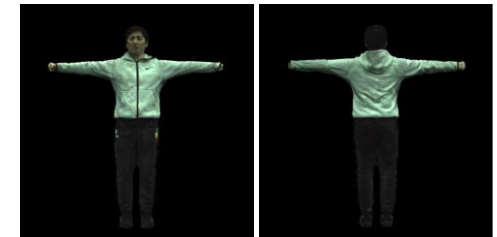
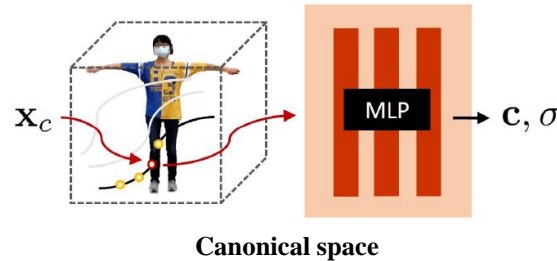
free-viewpoint rendering

Paper – HumanNeRF [CVPR2022]

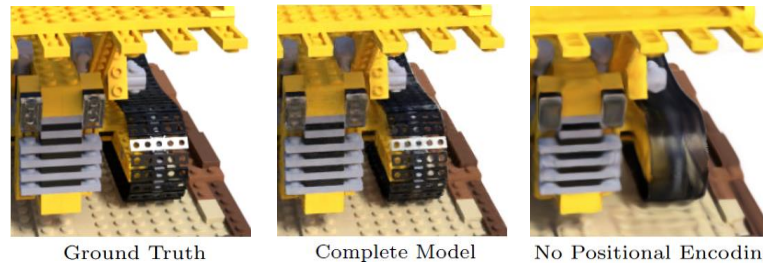
- 3D canonical space

- Canonical volume F_c

- Continuous field with an MLP
 - Outputs color \mathbf{c} and density σ
 - $F_c(\mathbf{x}) = \text{MLP}_{\theta_c}(\gamma(\mathbf{x}))$



Optimized canonical appearance
[ZJU-MoCaP-Subject387]



Positional Encoding

- Positional encoding $\gamma(\mathbf{x})$

- From NeRF¹), basic implementation is inefficient in the required number of samples per ray.
 - Positional encoding maps each input 5D coordinate into a higher dimensional space
 - It enables the MLP to represent higher frequency functions

⚡ Sinusoidal embedding function

$$\mathbf{x} \mapsto \{\sin(2^0\pi\mathbf{x}), \cos(2^0\pi\mathbf{x}), \dots, \sin(2^{L-1}\pi\mathbf{x}), \cos(2^{L-1}\pi\mathbf{x})\}$$

Paper – HumanNeRF [CVPR2022]

- Motion fields

- Transformation between observation field and canonical space
- To handle complex human movement with complex deformation by decomposing the motion field two parts

- Skeletal motion field T_{skel} : Inverse (volumetric) linear-blend skinning
- Non-rigid motion field T_{NR} : Complex deformation of non-rigid human appearance

- Skeletal motion field T_{skel} ¹⁾

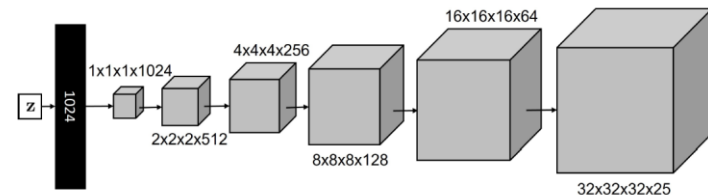
- Skeletal motion field provides the coarse deformation driven by standard skinning.

- $T_{skel}(\mathbf{x}, \mathbf{p}) = \sum_{i=1}^K \omega_o^i(\mathbf{x})(R_i \mathbf{x} + \mathbf{t}_i)$

- $\omega_o^i(\mathbf{x}) = \frac{\omega_c^i(R_i \mathbf{x} + \mathbf{t}_i)}{\sum_{i=1}^K \omega_c^i(R_i \mathbf{x} + \mathbf{t}_i)}$

∴ Solving for a single set of weight volumes $\{\omega_c^i(\mathbf{x})\}$ in canonical space can lead better generalization as it avoids over-fitting.

- $\{\omega_c^i(\mathbf{x})\} := W_c(\mathbf{x}) = \text{CNN}_{\theta_{skel}}(\mathbf{x}; \mathbf{z})$



Paper – HumanNeRF [CVPR2022]

• Motion fields

▪ Non-rigid motion field T_{NR}

– Output an offset Δx to the skeletal motion

$$-\Delta x(x, \mathbf{p}) = T_{NR}(T_{skel}(x, \mathbf{p}), \mathbf{p})$$

$$-T_{NR}(x, \mathbf{p}) = \text{MLP}_{\theta_{NR}}(\gamma(x); \Omega)$$

– Estimate difficult deformation of high-deformable region in dynamic human motion

▪ Delayed optimization of non-rigid motion field¹⁾

– Disable non-rigid motions at the beginning of optimization until 100K (in practice)

– Hann window was applied to frequency bands of positional encoding \rightarrow prevent over-fitting

$$\omega(\alpha) = \frac{1 - \cos(\text{clamp}(\alpha - j, 0, 1)\pi)}{2}, \quad \alpha(t) = L \frac{\max(0, t - T_s)}{T_e - T_s}, \quad j \in \{0, \dots, L - 1\}$$

⚡ $\omega(\alpha)$: weight for each frequency band j of positional encoding

⚡ $\alpha(t)$: width of a truncated Hann window

$$\gamma_\alpha(x) = \{w_0 \sin(2^0 \pi x), w_0 \cos(2^0 \pi x), \dots, w_{L-1} \sin(2^{L-1} \pi x), w_{L-1} \cos(2^{L-1} \pi x)\}$$

⚡ If $\alpha = 0$, non-rigid motion completely be disabled

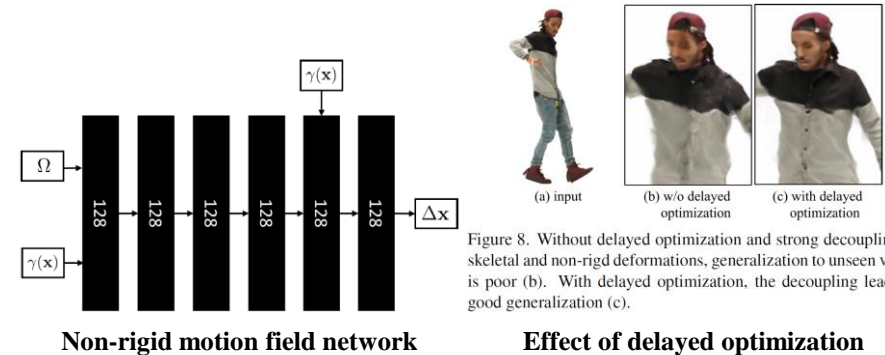


Figure 8. Without delayed optimization and strong decoupling of skeletal and non-rigid deformations, generalization to unseen views is poor (b). With delayed optimization, the decoupling leads to good generalization (c).

Paper – HumanNeRF [CVPR2022]

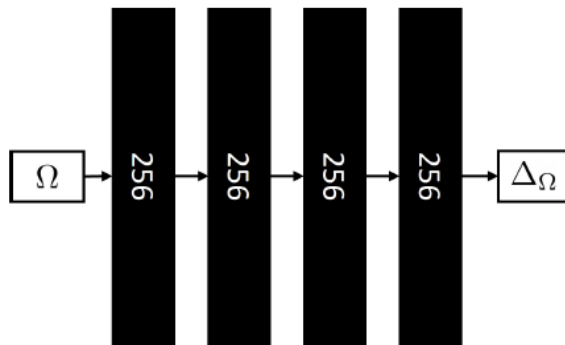
- Pose refinement module

- Pose refinement

- Initial human poses from ‘pose detector’ is not accuracy.
 - MLP based network outputs difference of joint rotative vectors $\Delta\Omega$.
 - Network better explains the observations and improve the skeleton-driven deformation.
 - Pose correction function P_{pose}

- ⌘ Consider 23 joints except for the root. (body orientation)

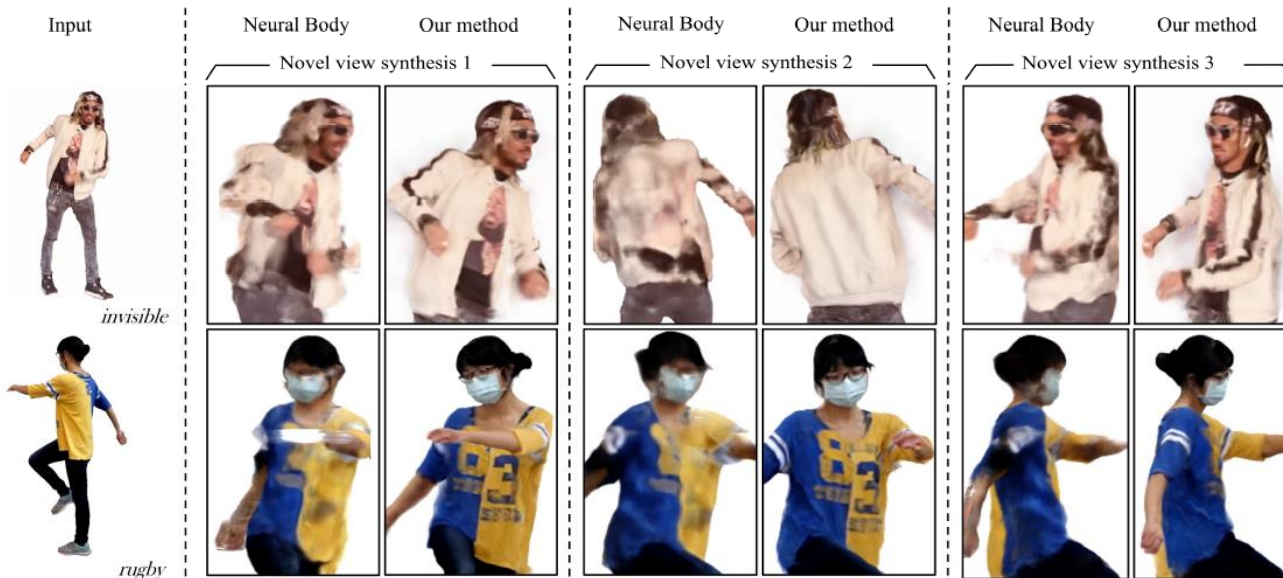
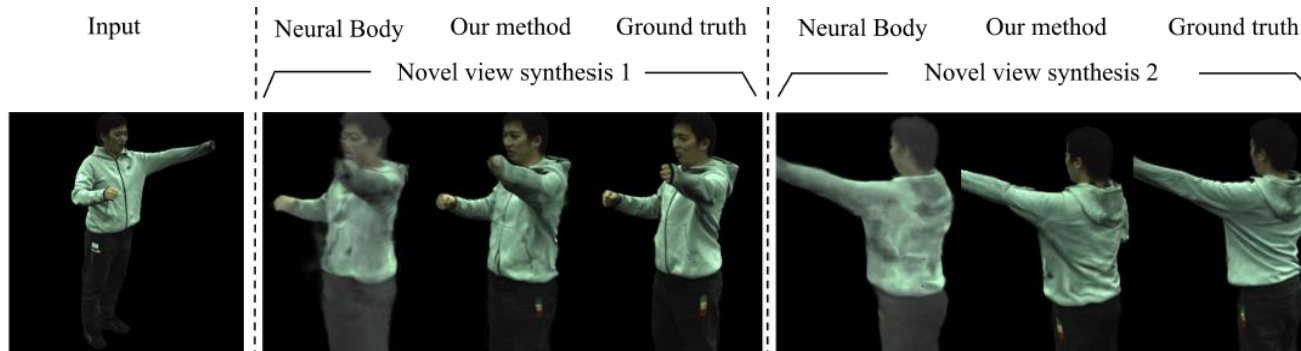
- ⌘ Instead, describe changes of global body orientation as camera rotations.



Pose refinement network architecture

Paper – HumanNeRF [CVPR2022]

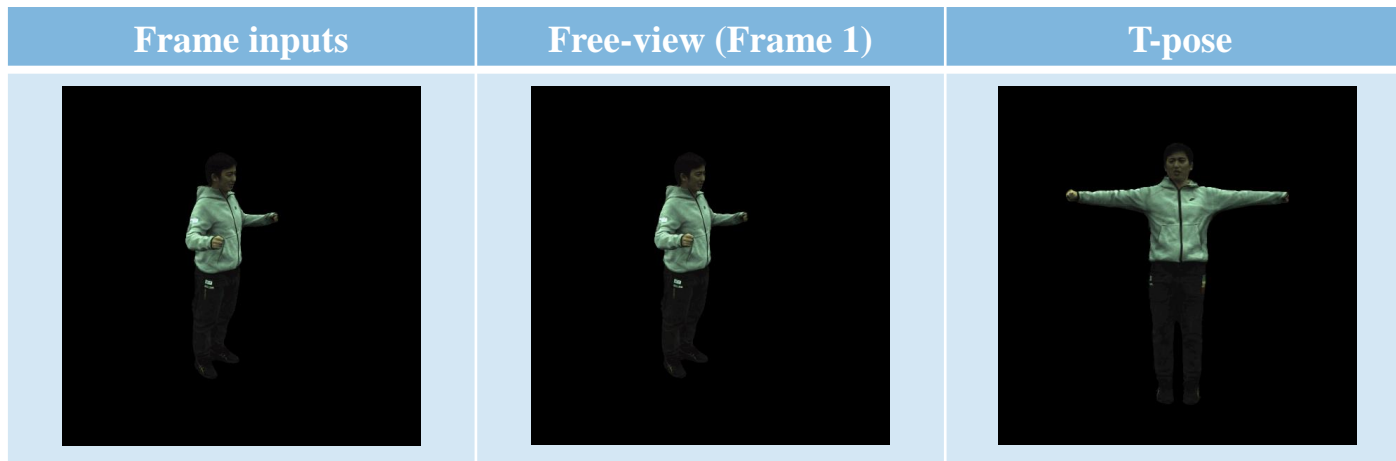
- Qualitative results



Paper – HumanNeRF [CVPR2022]

- Renderings

- Zju-MoCaP - subject 387



- Performance

	Subject 377			Subject 386			Subject 387		
	PSNR ↑	SSIM ↑	LPIPS* ↓	PSNR ↑	SSIM ↑	LPIPS* ↓	PSNR ↑	SSIM ↑	LPIPS* ↓
Neural Body [48]	29.11	0.9674	40.95	30.54	0.9678	46.43	27.00	0.9518	59.47
Ours	30.41	0.9743	24.06	33.20	0.9752	28.99	28.18	0.9632	35.58

	Subject 392			Subject 393			Subject 394		
	PSNR ↑	SSIM ↑	LPIPS* ↓	PSNR ↑	SSIM ↑	LPIPS* ↓	PSNR ↑	SSIM ↑	LPIPS* ↓
Neural Body [48]	30.10	0.9642	53.27	28.61	0.9590	59.05	29.10	0.9593	54.55
Ours	31.04	0.9705	32.12	28.31	0.9603	36.72	30.31	0.9642	32.89

Paper – HumanNeRF [CVPR2022]

- Limitations

- Occlusion

- It has artifacts when part of the body is not shown in the video (occlusion)

- Pose correction failure

- If the initial pose estimate is poor or the image contains strong artifact such as motion blur

- Non-rigid motion covering

- Authors assume non-rigid motion is pose-dependent, but it is not always true

- E.g. clothes shifting due to wind or due to follow-through after dynamic subject motion

- Only consider human scene, not background scene