

Virtual Humans – Winter 23/24

Lecture 10_2 – Humans and NeRF

Prof. Dr.-Ing. Gerard Pons-Moll

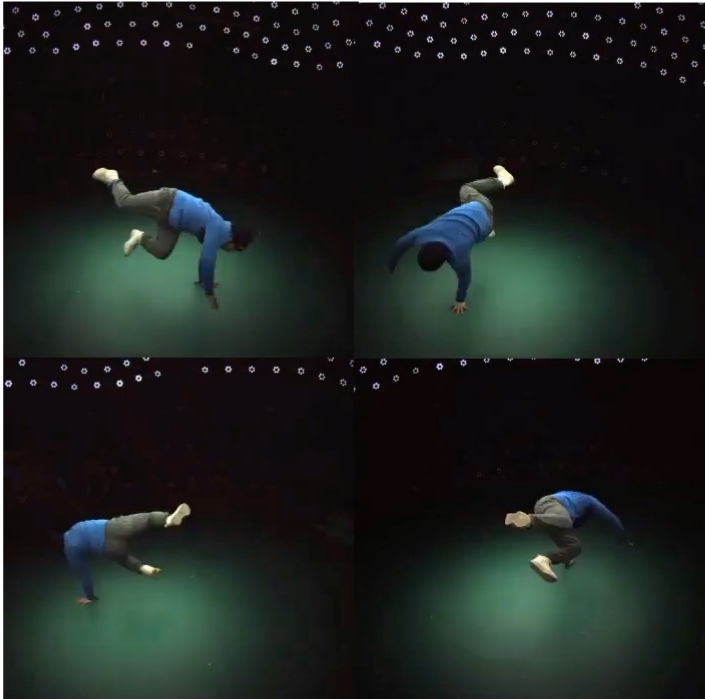
University of Tübingen / MPI-Informatics

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN



Novel View Synthesis for Humans

Task: Novel view synthesis from a sparse multi-view video



4-view video



Novel view synthesis of dynamic human
(Our result)

Peng et al. CVPR 21

Human models using NeRF

Challenge: It is ill-posed to learn 3D representations from very sparse observations



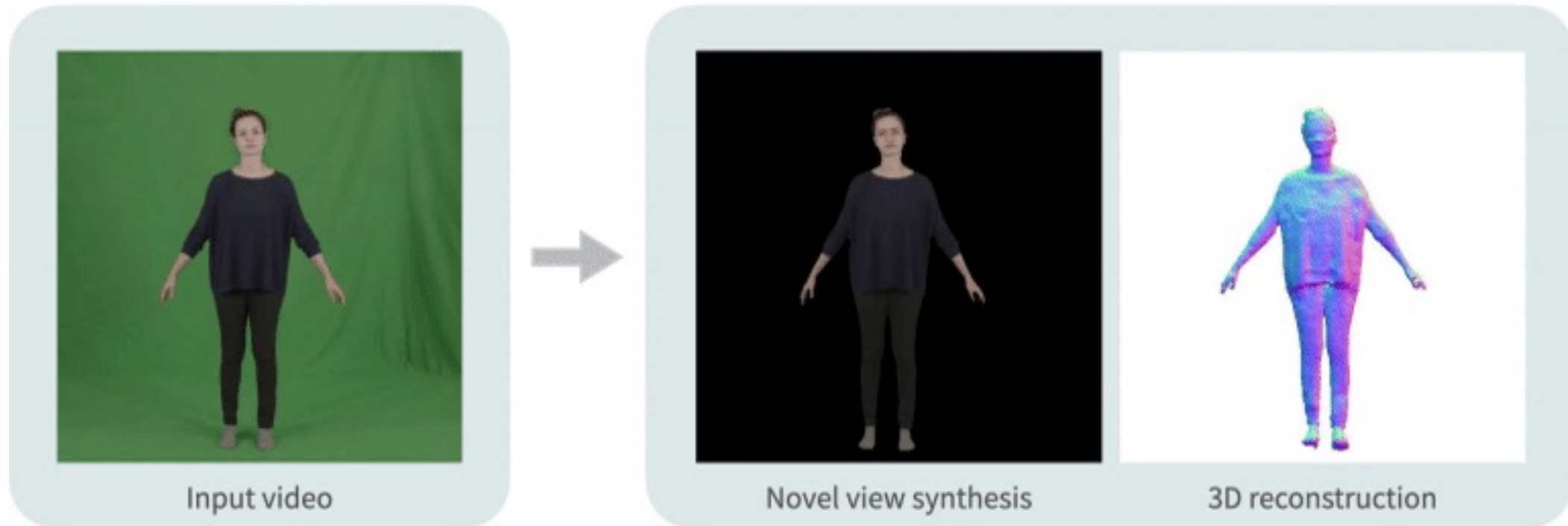
Four input images



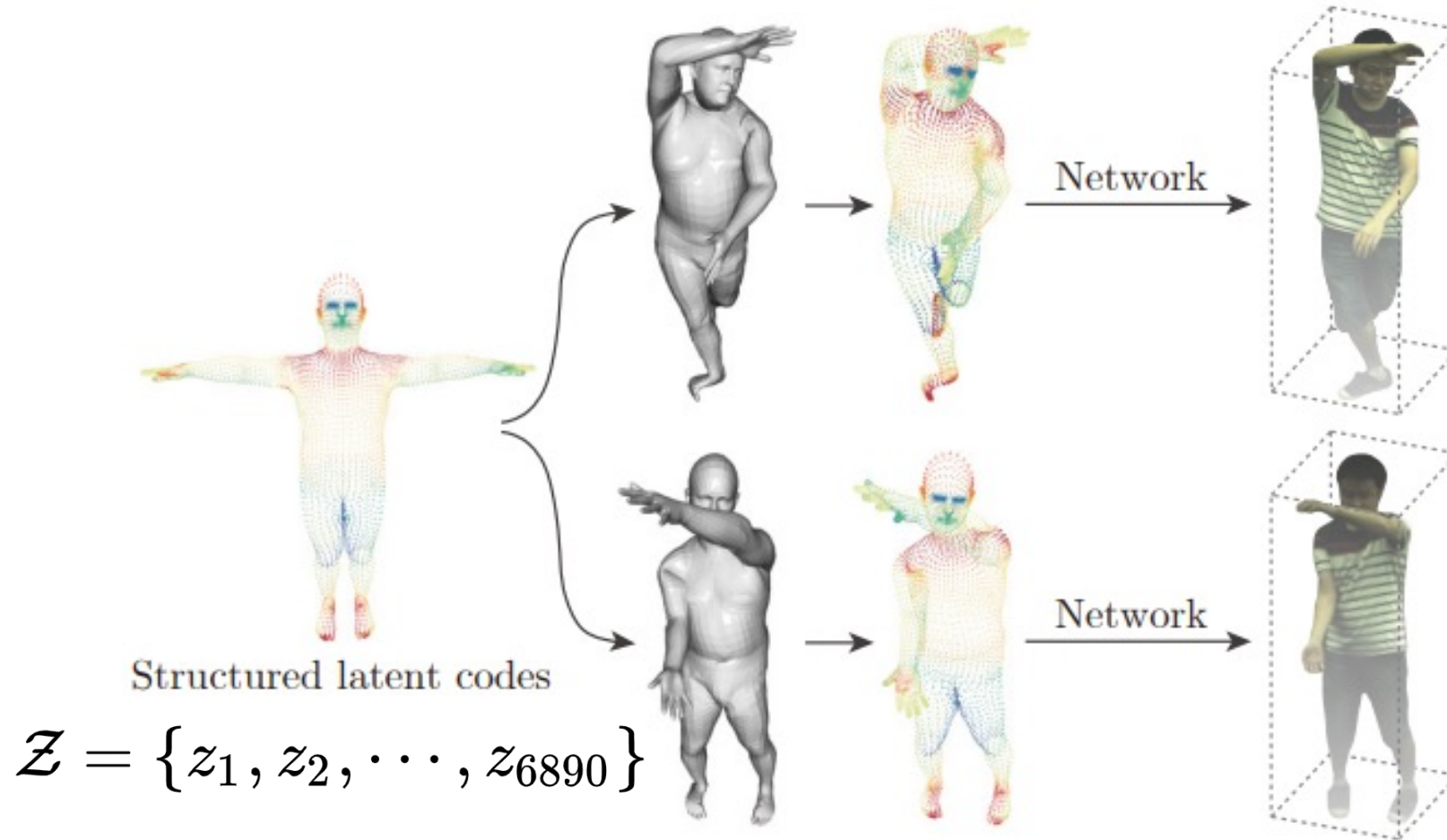
Novel view synthesis by NeRF [3]

[3] Mildenhall, Ben, et al. Nerf: Representing scenes as neural radiance fields for view synthesis. In ECCV, 2020.

Neural Body: Implicit Neural Representations with Structured Latent Codes for NVS of Dynamic Humans

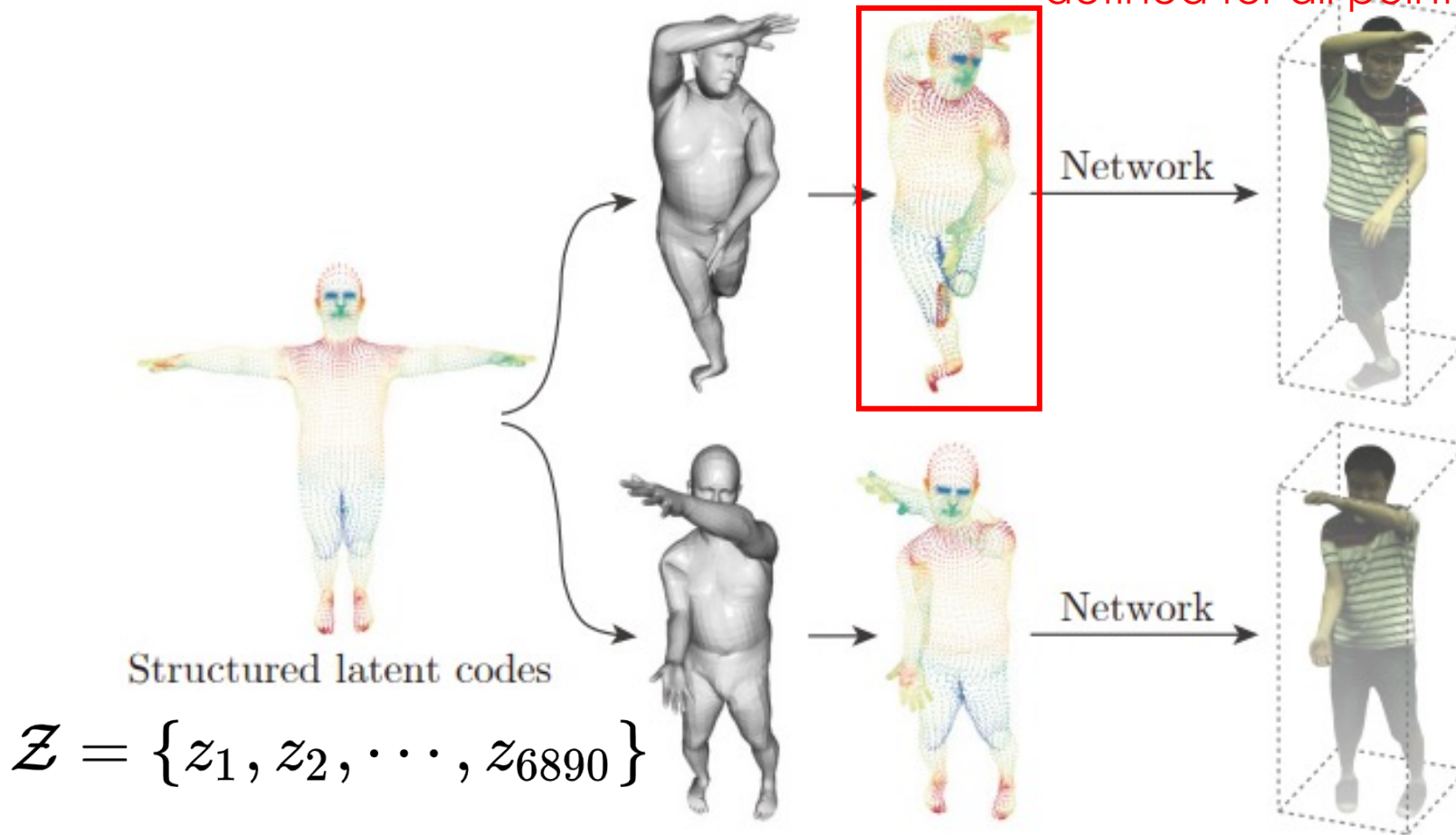


Neural Body: Key Idea

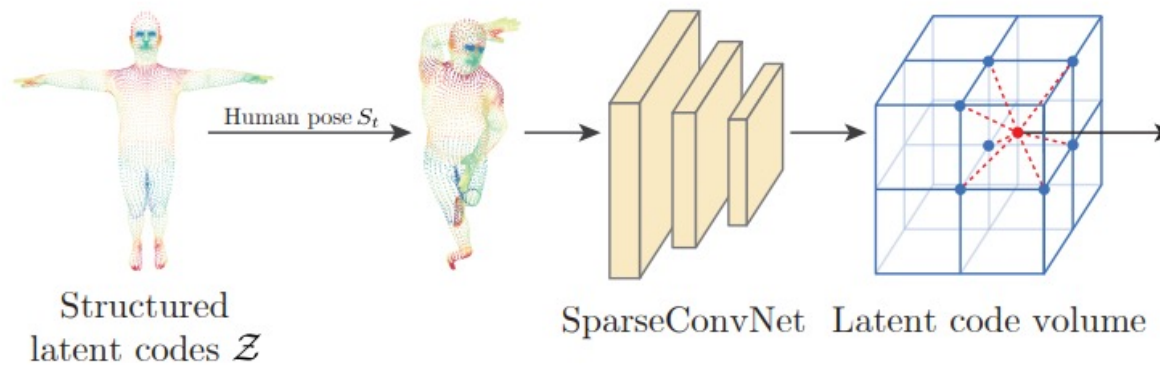


Neural Body: Key Idea

Sparse latent code, not defined for all points in 3D

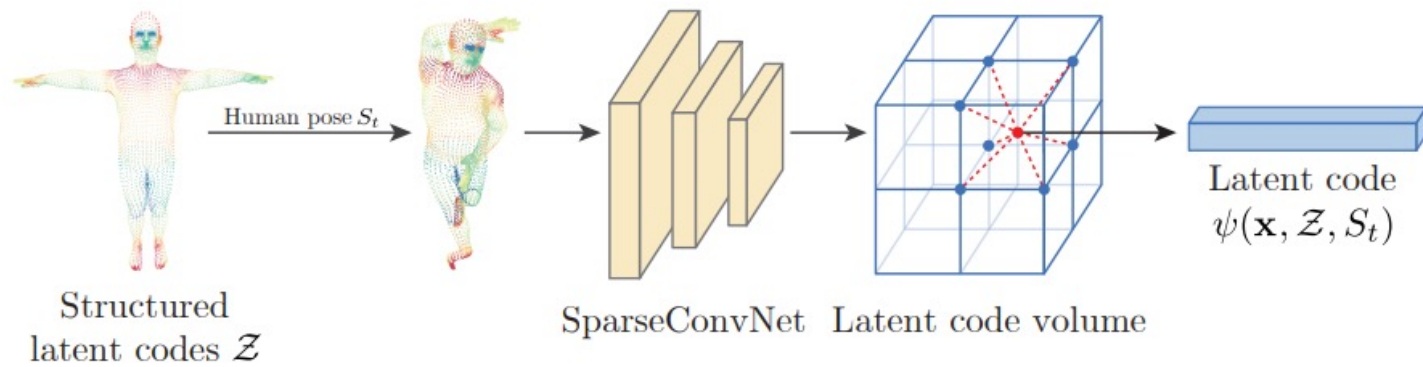


Neural Body: Method



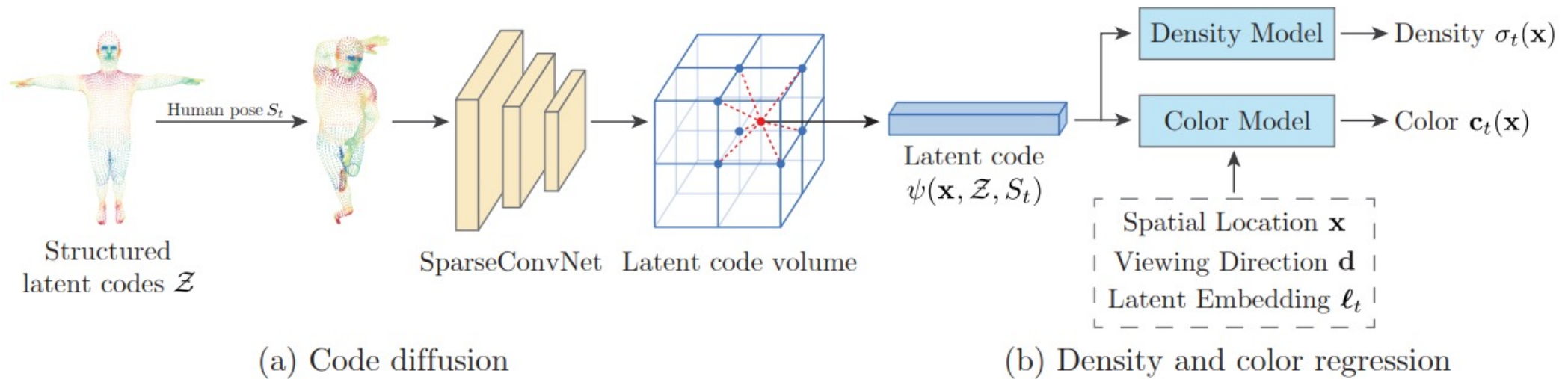
(a) Code diffusion

Neural Body: Method



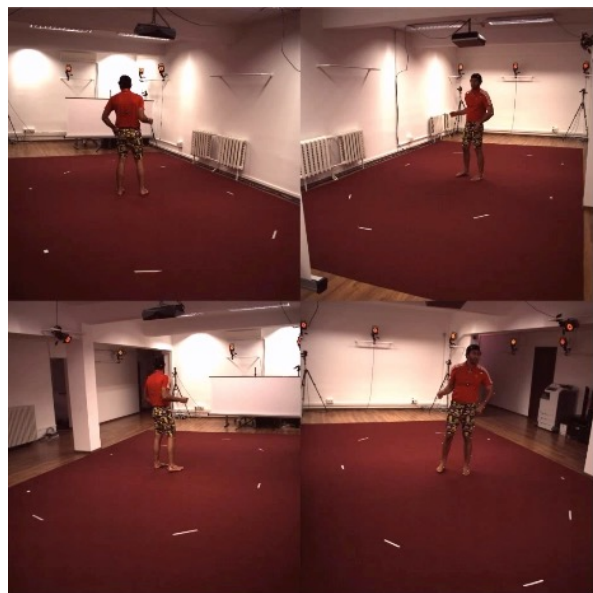
(a) Code diffusion

Neural Body: Method



Neural Body: Results

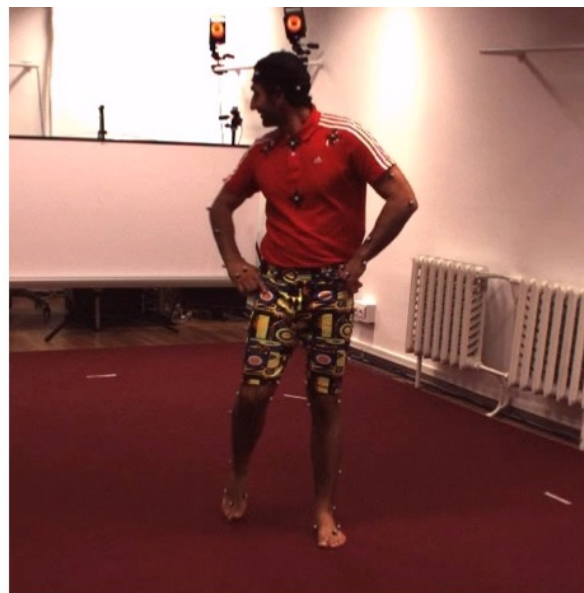
Input Views



NeuralBody



3D reconstruction



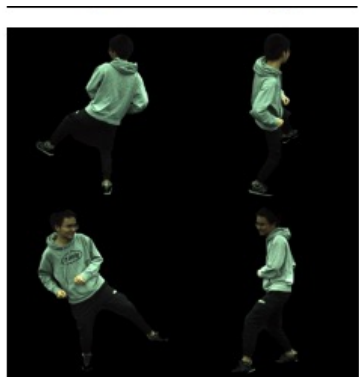
Neural Body: Conclusion

- Use SMPL mesh as structure:
 - + Strong human prior and preserves human shape.



Neural Body: Conclusion

- Use SMPL mesh as structure:
 - + Strong human prior and preserves human shape.
 - Introduces artifacts in clothing and complex motions that are not captured by the SMPL model.



Input views



OURS



PIFuHD

Input



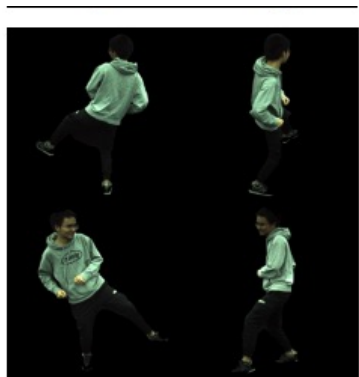
invisible

NeuralBody



Neural Body: Conclusion

- Use SMPL mesh as structure:
 - + Strong human prior and preserves human shape.
 - Introduces artifacts in clothing and complex motions that are not captured by the SMPL model.
 - Only works for multi-view setup.



Input views



OURS

PIFuHD

Input

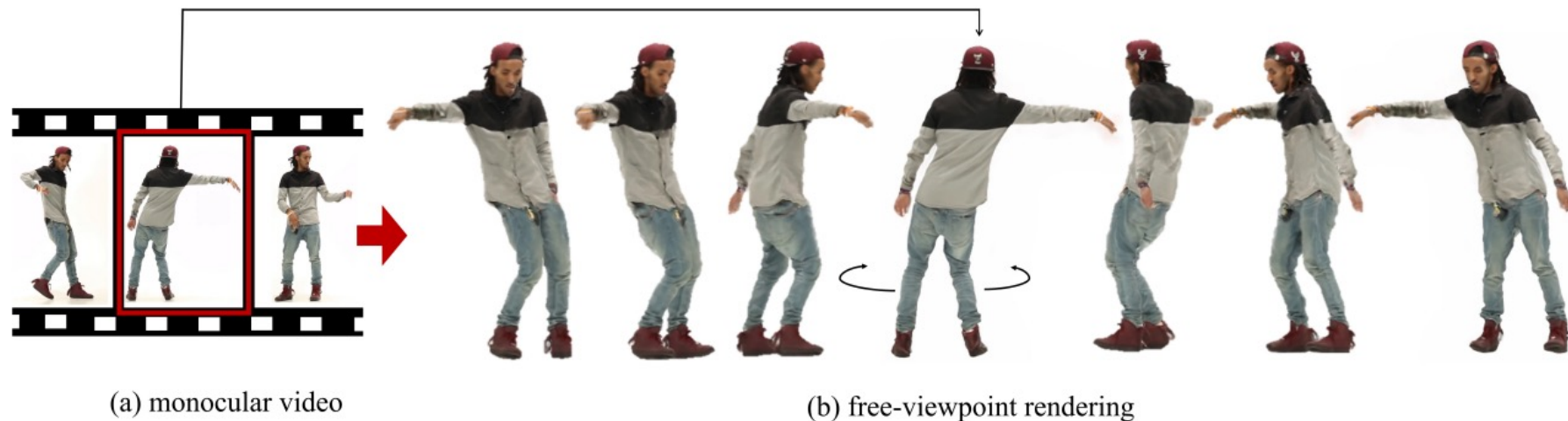


invisible

NeuralBody



HumanNeRF: Free-viewpoint Rendering of Moving People from Monocular Video



- Given a monocular video (a) of a human performing complex movement, e.g., dancing (left), HumanNeRF creates a free-viewpoint rendering for any frame in the sequence (b).
- Deformation and skinning formulation similar to NeuralGIF

HumanNeRF: Key Idea

- Split the deformation into:
 1. Human articulation
 2. Non-rigid pose dependent deformation

Similar to
NeuralGIF

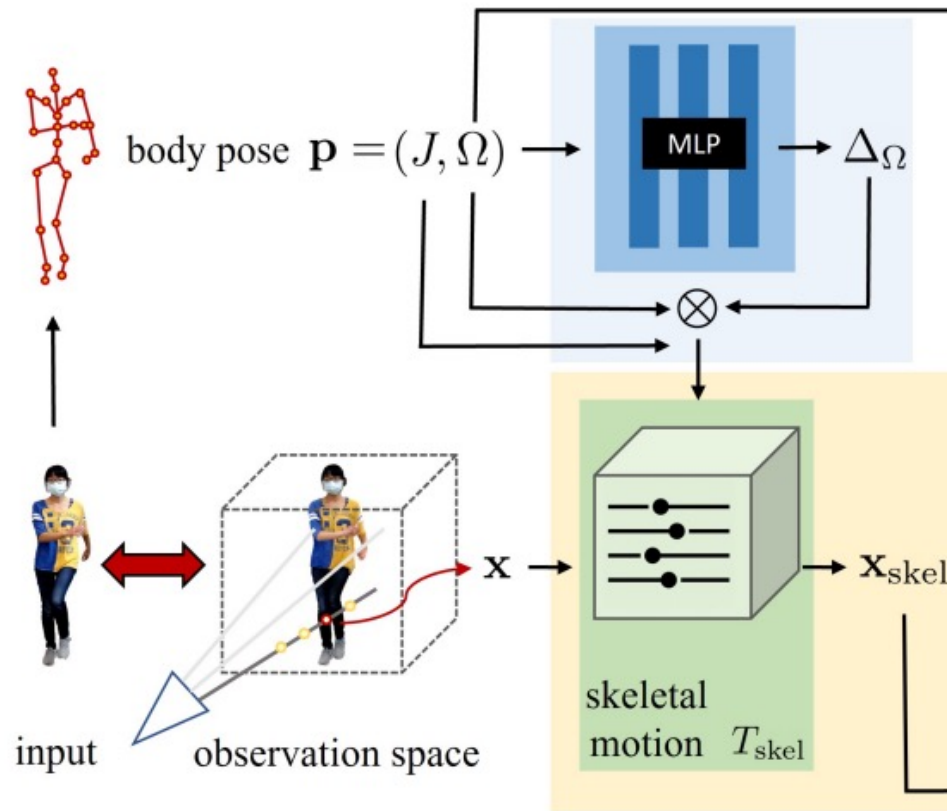
HumanNeRF: Key Idea

- Split the deformation into:
 1. Human articulation
 2. Non-rigid pose dependent deformation
- Skinning weights using forward skinning.

Similar to
NerualGIF

Similar to SNARF

HumanNeRF: Method

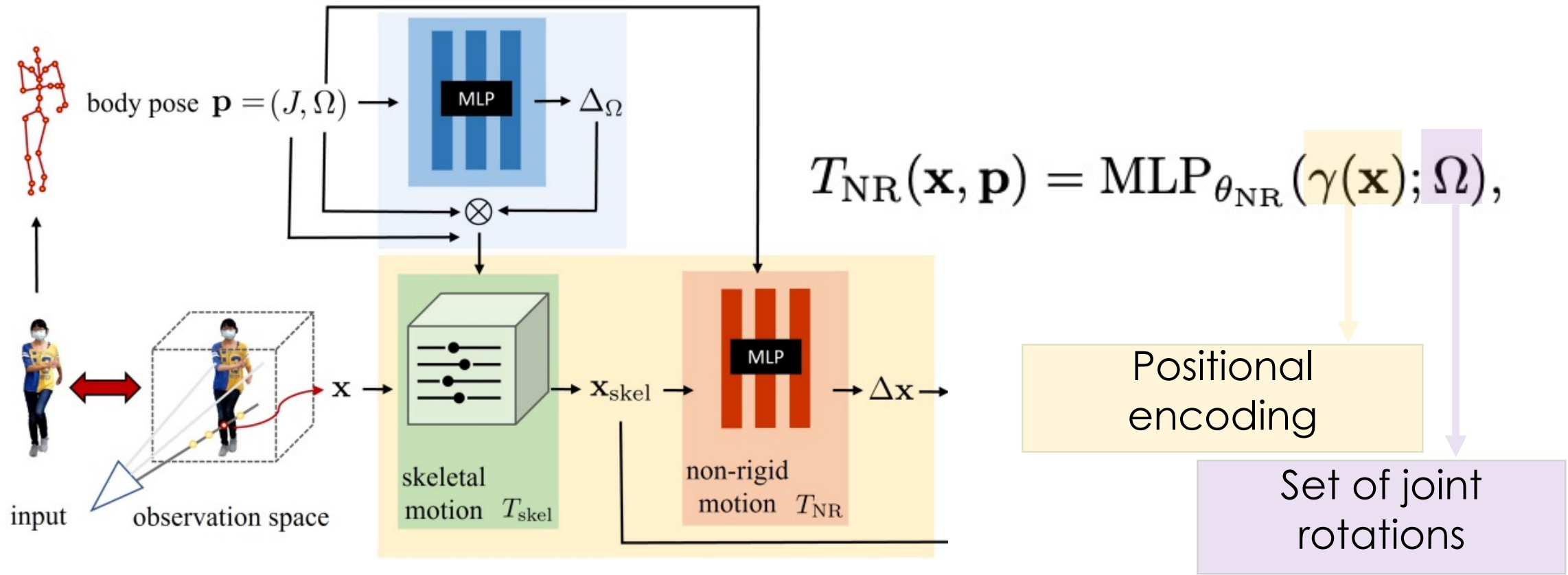


$$T_{\text{skel}}(\mathbf{x}, \mathbf{p}) = \sum_{i=1}^K w_o^i(\mathbf{x}) (R_i \mathbf{x} + \mathbf{t}_i),$$

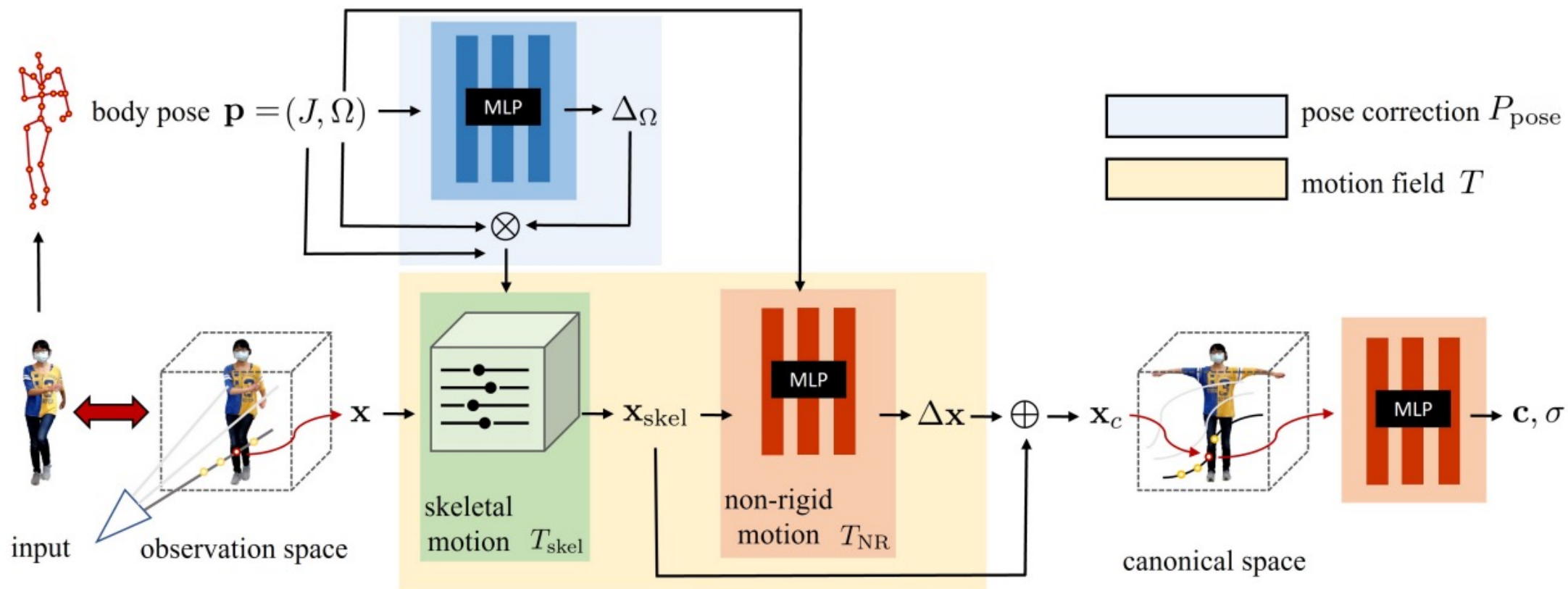
Observed to
canonical

Skinning weight obtained
using forward skinning

HumanNeRF: Method



HumanNeRF: Method



HumanNeRF: Method



More on Human and NeRFs

- Animatable NeRF, Peng et al., ICCV2021
- H-NeRF, Xu et al., NeurIPS 2021
- NeuMan, Jian et al., ECCV 2022
- DoubleField, Shao et al., CVPR 2022
-
- And many more.

Limitation of NeRF/Implicit Representations

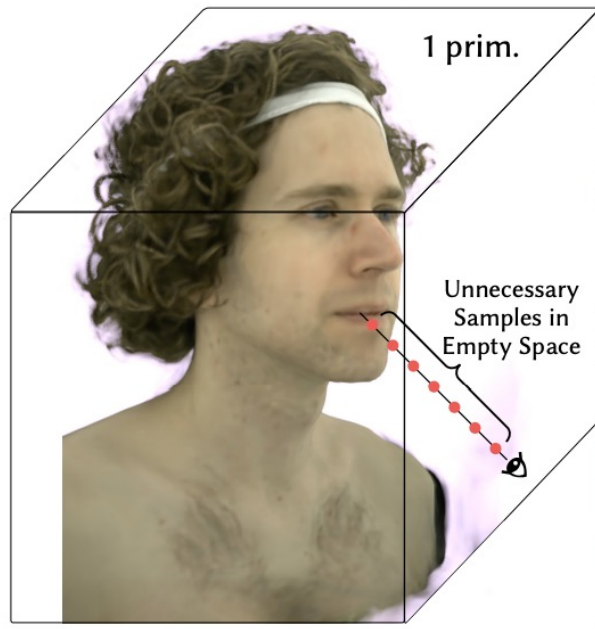
- **4. Expensive training:**
 - Training is slow(10 hours-upto few days)
 - Inference is also not real time

Mixture of Volumetric Primitives

- Combines the advantages of volumetric and primitive-based approaches for:
 - High performance decoding
 - Efficient rendering
- A novel motion model for voxel grids for scene motion, minimization of primitive overlap to increase the representational power.

Mixture of Volumetric Primitives: Key Idea

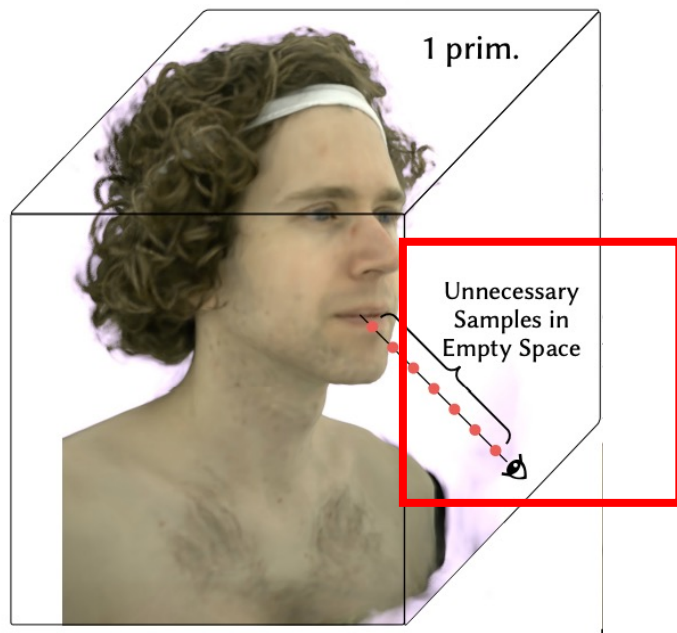
- Combining primitives and volumetric representation:



Single Volume
2M voxels

Mixture of Volumetric Primitives: Key Idea

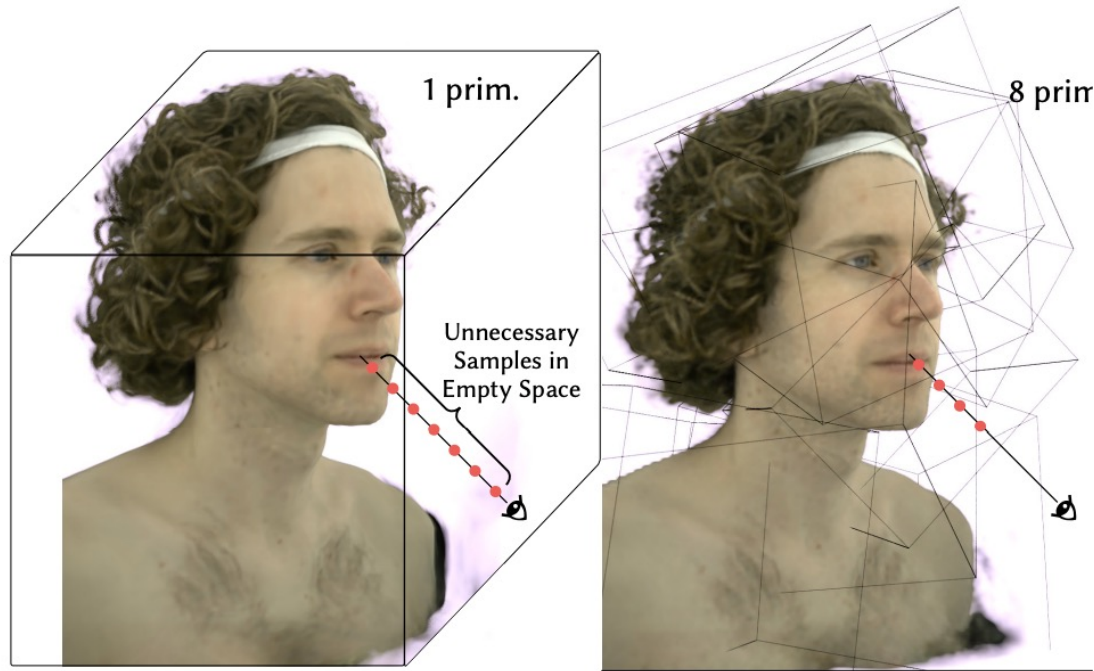
- Combining primitives and volumetric representation:



Single Volume
2M voxels

Mixture of Volumetric Primitives: Key Idea

- Combining primitives and volumetric representation:



Single Volume
2M voxels

Mixtures of Volumetric Primitives
2M voxels

$$\mathcal{V}_k = \{t_k, R_k, s_k, V_k\}$$

Translation, rotation
and scale of
grid/primitives

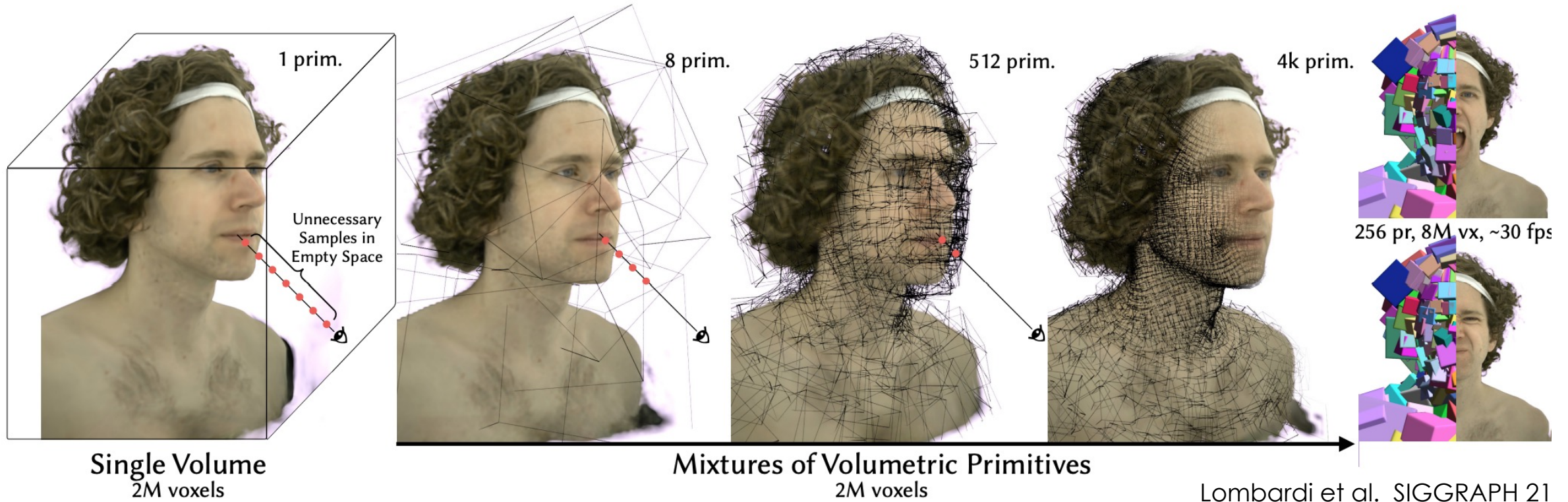
$$V_k \in \mathbb{R}^{4 \times M_x \times M_y \times M_z}$$

Where each feature grid
contains color and density
and M is grid resolution

Lombardi et al. SIGGRAPH 21

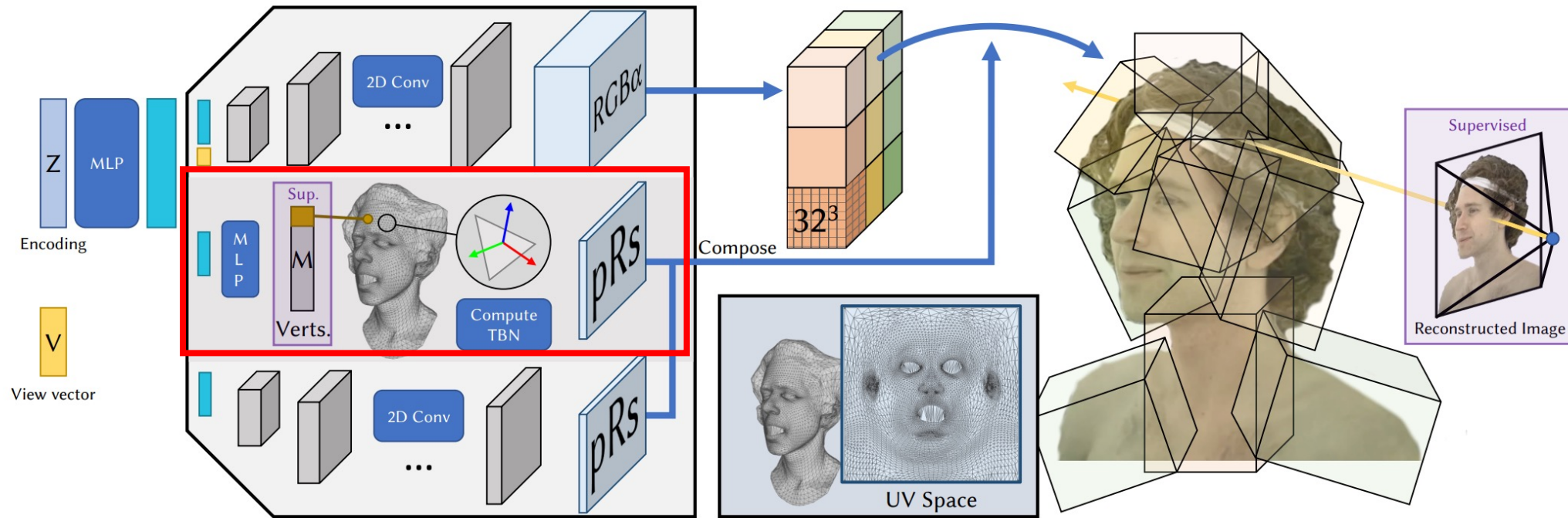
Mixture of Volumetric Primitives: Key Idea

- Combining primitives and volumetric representation:

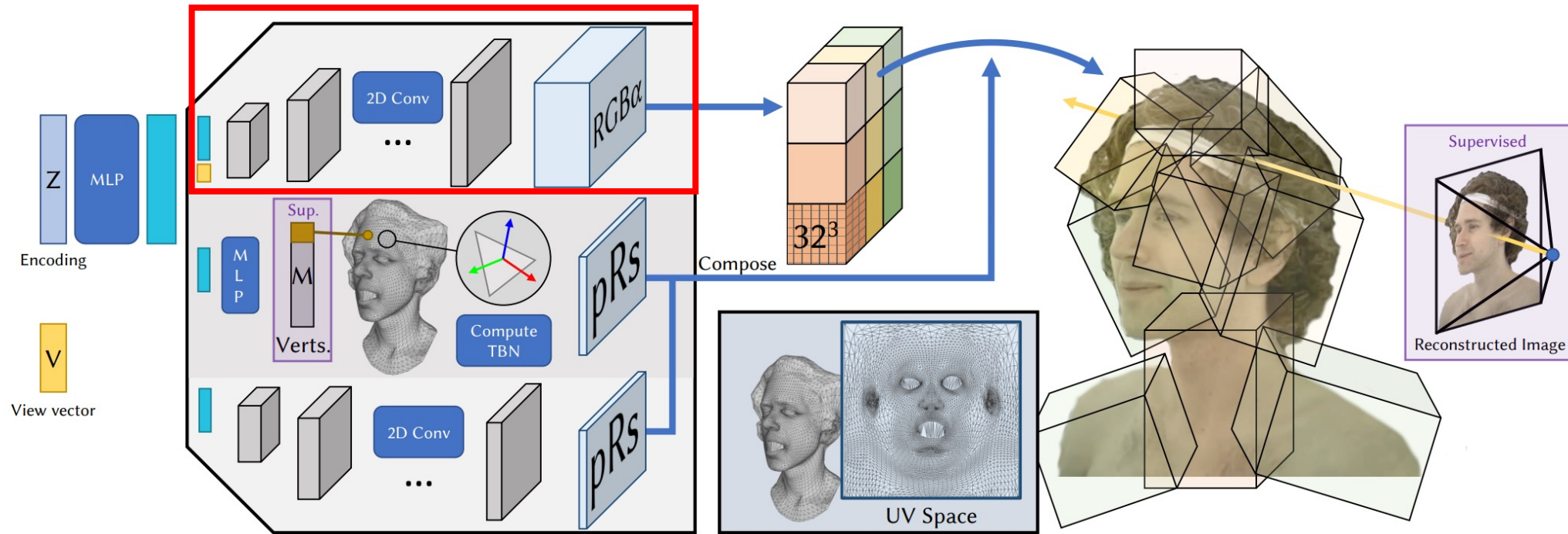


Mixture of Volumetric Primitives: Method

Mixture of Volumetric Primitives: Method



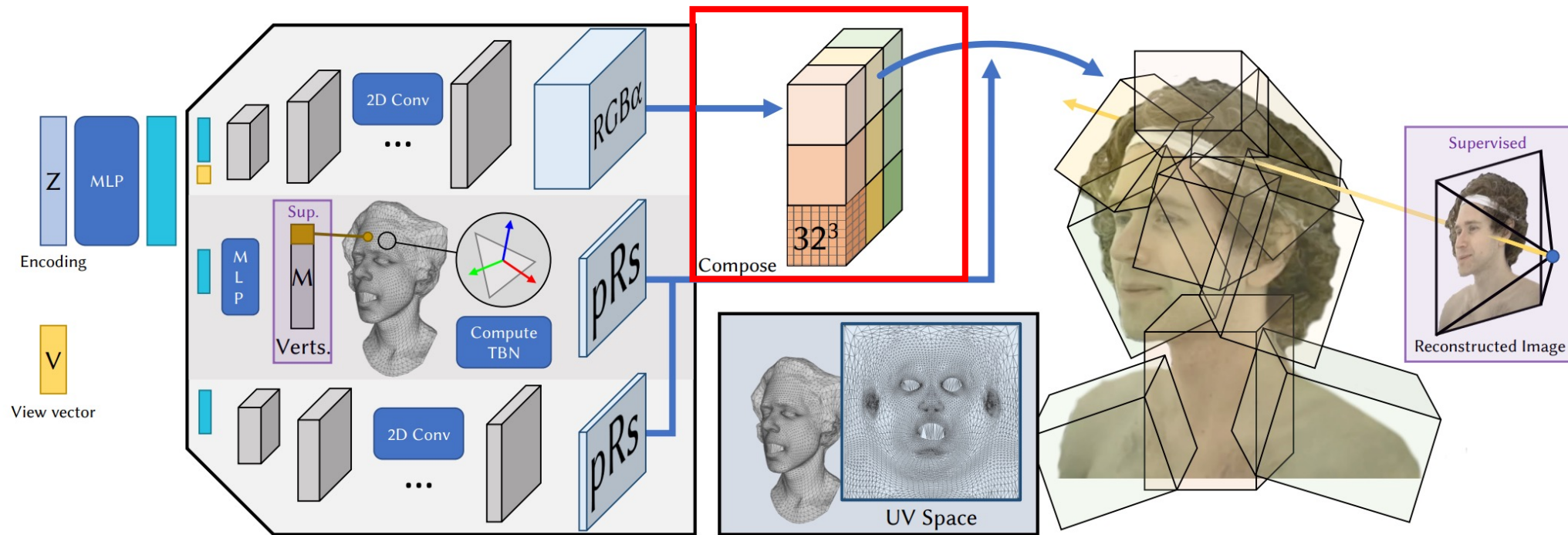
Mixture of Volumetric Primitives: Method



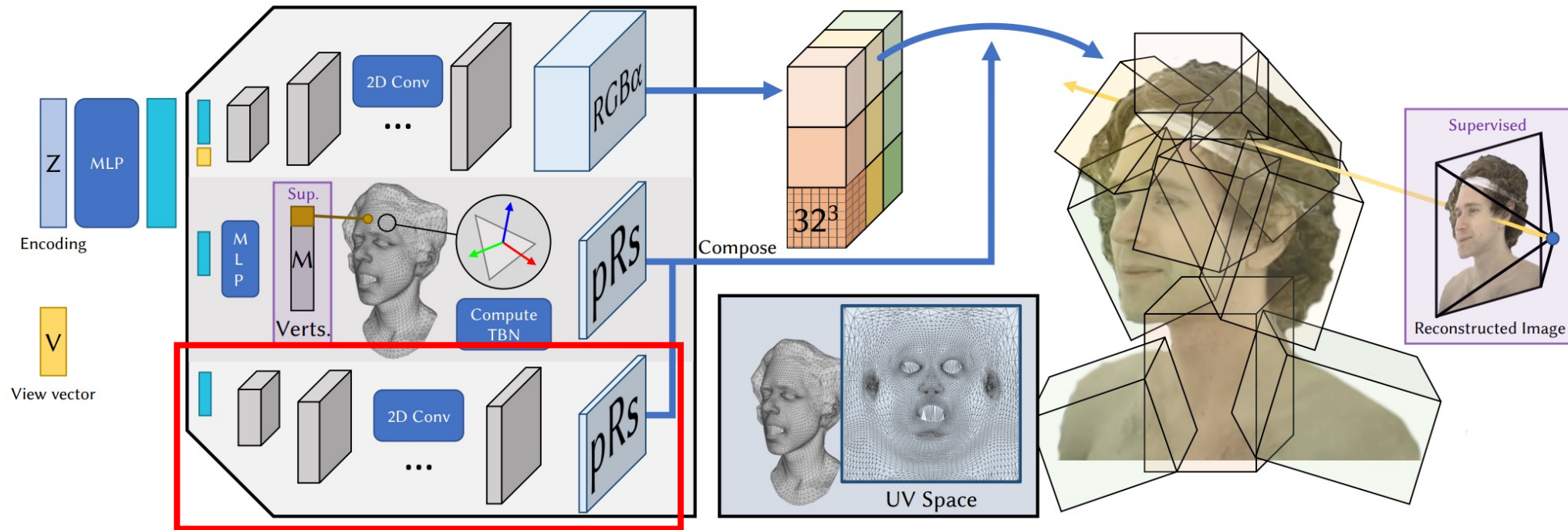
$$V_k \in \mathbb{R}^{4 \times M_x \times M_y \times M_z}$$

Where each feature grid contains color and density and M is grid resolution

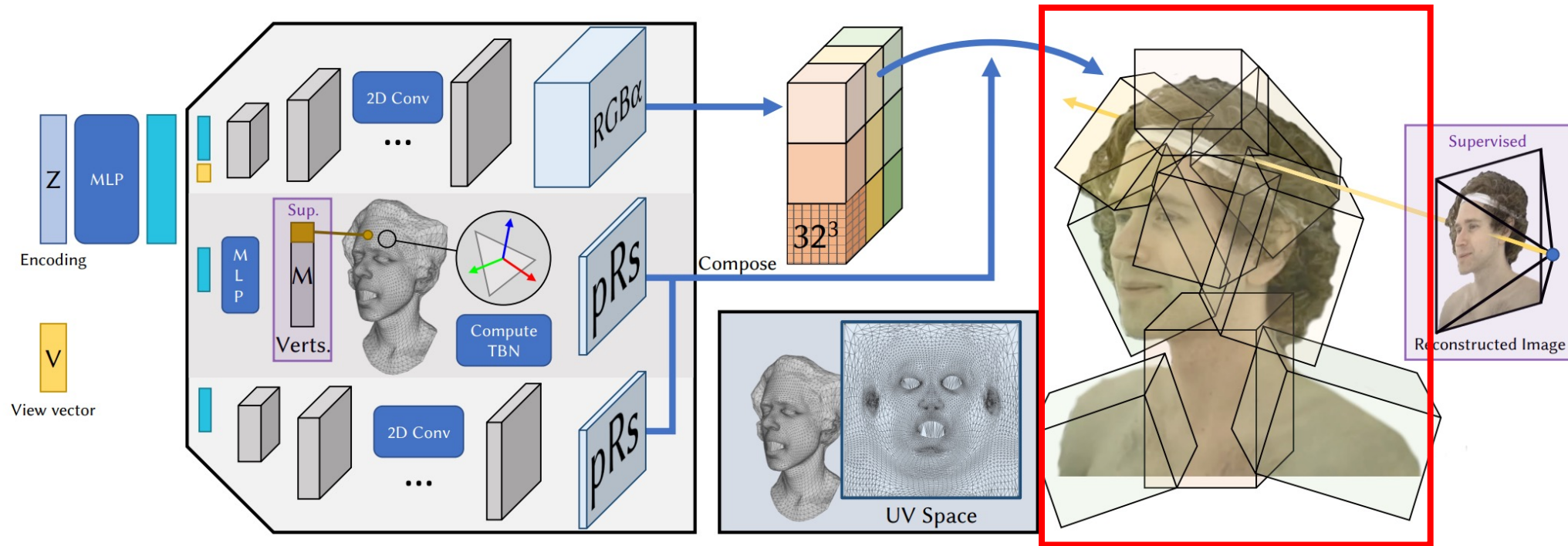
Mixture of Volumetric Primitives: Method



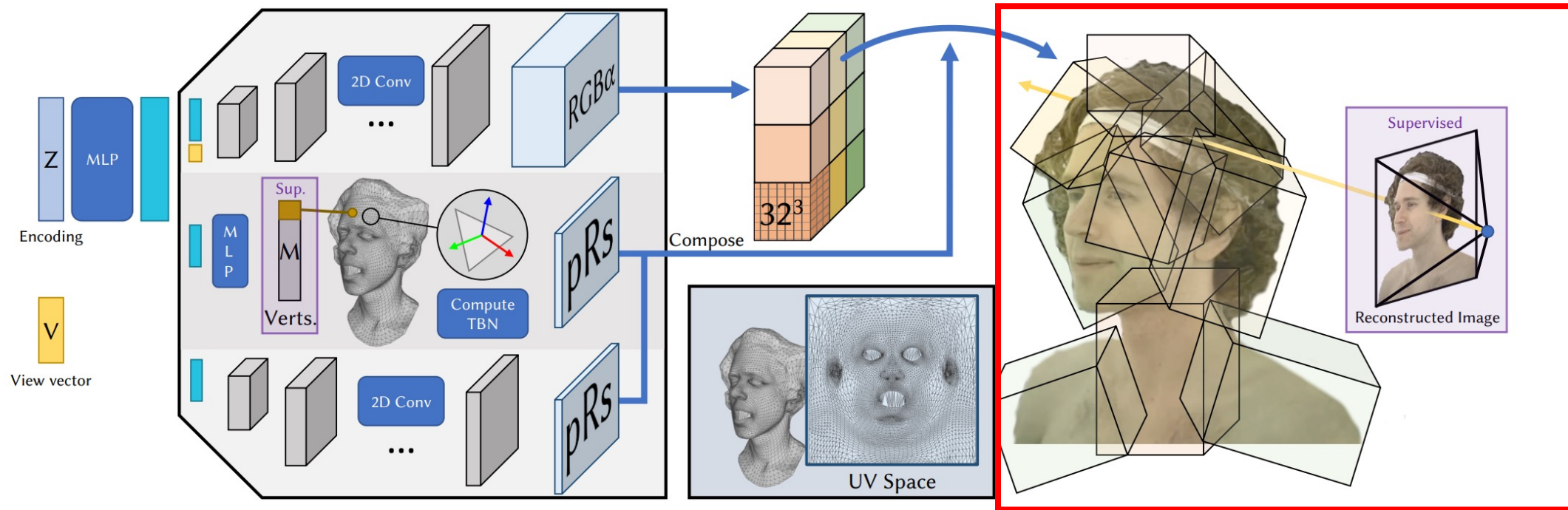
Mixture of Volumetric Primitives: Method



Mixture of Volumetric Primitives: Method



Mixture of Volumetric Primitives: Method



Mixture of Volumetric Primitives: Training

$$\mathcal{L}(\Theta; \mathcal{I}_p) = \mathcal{L}_{\text{pho}}(\Theta; \mathcal{I}_p) + \mathcal{L}_{\text{geo}}(\Theta) + \mathcal{L}_{\text{vol}}(\Theta) + \mathcal{L}_{\text{del}}(\Theta) + \mathcal{L}_{\text{kld}}(\Theta)$$

$$\mathcal{L}_{\text{pho}} = \lambda_{\text{pho}} \frac{1}{N_{\mathcal{P}}} \sum_{p \in \mathcal{P}} \|\mathcal{I}_p - \bar{\mathcal{I}}_p(\Theta)\|_2^2$$

Difference between predicted and GT image

$$\mathcal{L}_{\text{vol}} = \lambda_{\text{vol}} \sum_{i=1}^{N_{\text{prim}}} \text{Prod}(\mathbf{s}_i)$$

Volumetric primitive to be as small as possible

$$\mathcal{L}_{\text{geo}} = \lambda_{\text{geo}} \frac{1}{N_{\text{mesh}}} \sum_{i=0}^{N_{\text{mesh}}} \|\mathbf{v}_i - \bar{\mathbf{v}}_i(\Theta)\|_2^2$$

Difference regressed vertex position and GT vertex

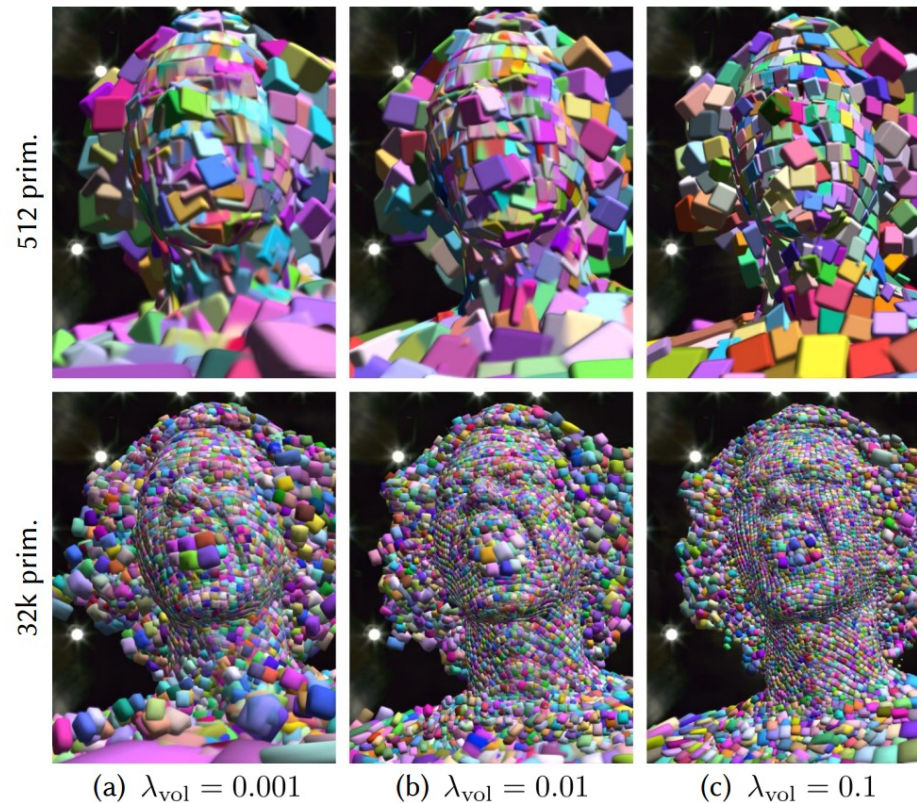
Mixture of Volumetric Primitives: Results



Lombardi et al. SIGGRAPH 21

Mixture of Volumetric Primitives: Results

A stronger primitive volume prior leads to less overlap and thus speeds up raymarching



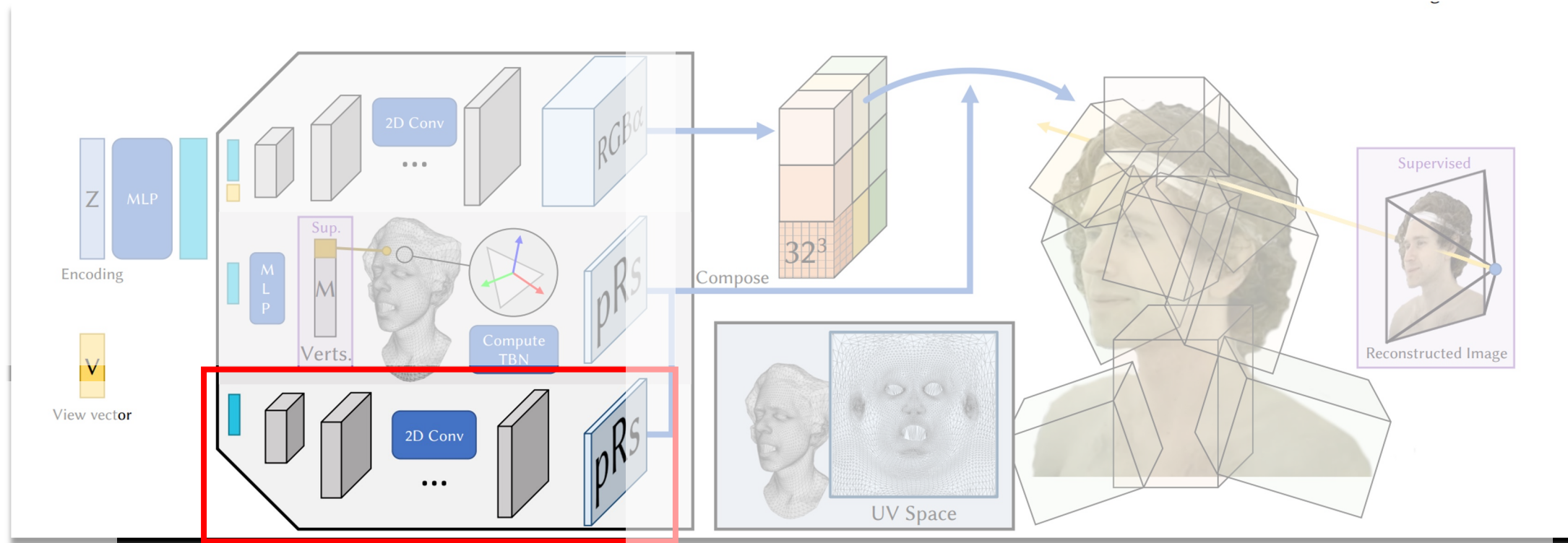
Lombardi et al. SIGGRAPH 21

Mixture of Volumetric Primitives

How much the result depends on initialization(or the guide coarse mesh)?

Mixture of Volumetric Primitives

How much the result depends on initialization(or the guide coarse mesh)?



Lombardi et al. SIGGRAPH 21

Mixture of Volumetric Primitives

How much the result depends on initialization(or the guide coarse mesh)?

- **Alpha Fade:** Windowing function adds an inductive bias to explain the scene's contents via motion instead of volumetric opacity.

$$W(x, y, z) = \exp\left(-\alpha(x^\beta + y^\beta + z^\beta)\right)$$

$$W(x, y, z) \in \mathbb{R}^{M^3}$$

Mixture of Volumetric Primitives

How much the result depends on initialization(or the guide coarse mesh)?



(a) GT



(b) No alpha fade



(c) With alpha fade

Lombardi et al. SIGGRAPH 21

Mixture of Volumetric Primitives

+ Combine volumetric and primitive based approach for generalizable representation of dynamic scenes.

- Fast to render
- Represent translucent parts, thing structures

Mixture of Volumetric Primitives

- + Combine volumetric and primitive based approach for generalizable representation of dynamic scenes.
 - Fast to render
 - Represent translucent parts, thing structures
- + Strong prior about structure of underlying shape via coarse shape and derived primitives

Mixture of Volumetric Primitives

+ Combine volumetric and primitive based approach for generalizable representation of dynamic scenes.

- Fast to render
- Represent translucent parts, thing structures

+ Strong prior about structure of underlying shape via coarse shape and derived primitives

- No prior/information about the structure of motion of underlying object.

- Difficulty to model human motion

Drivable Volumetric Avatars

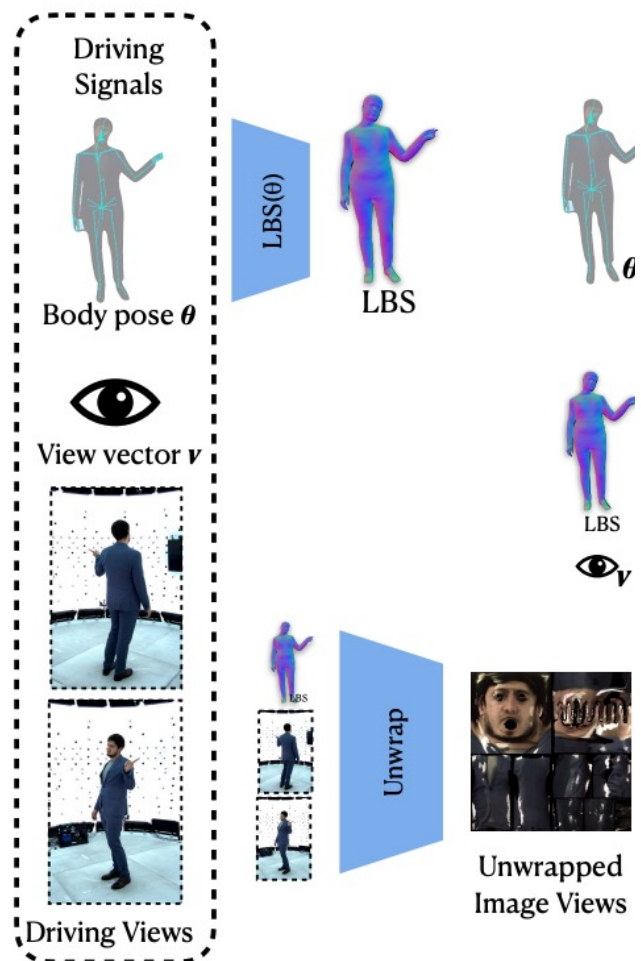


Drivable Volumetric Avatars

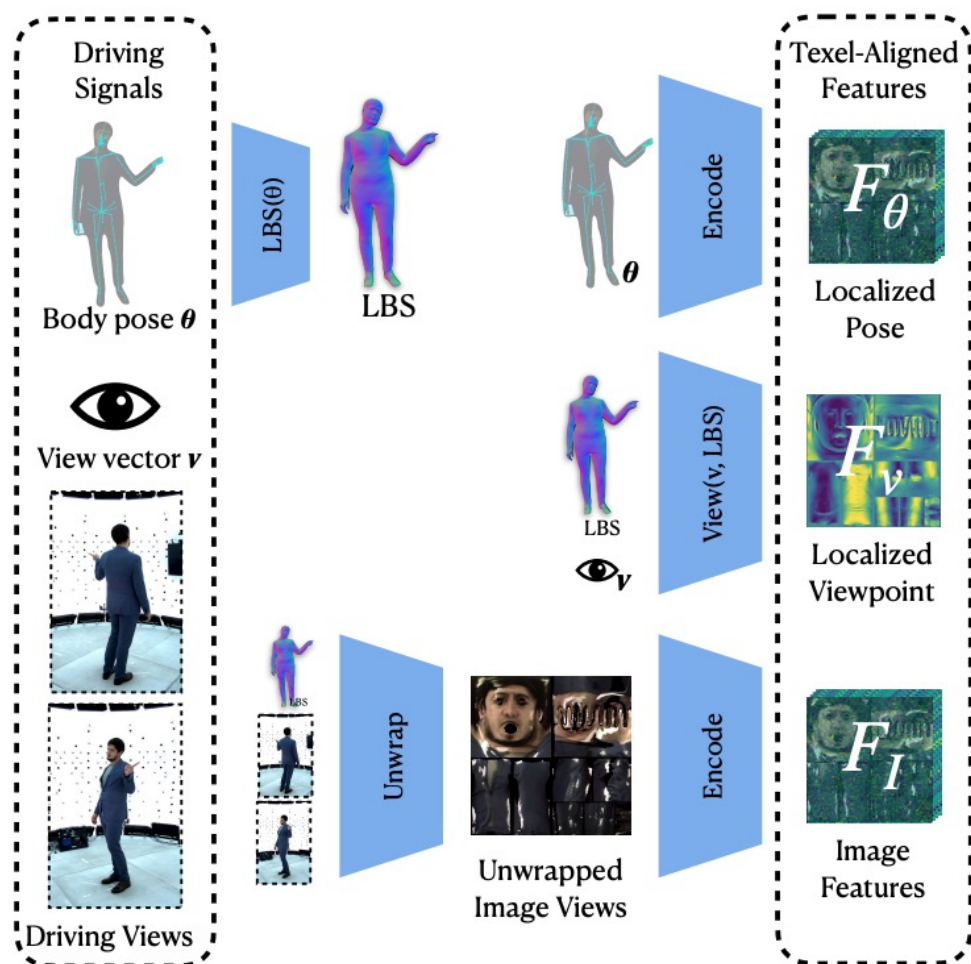
Key Idea: Articulated Primitives

- Use SMPL posed mesh. M_{θ} as a guide mesh and define primitives using it.
- Initialise primitive by uniformly sampling UV space and mapping each primitive to closest texel $\hat{t}_k(\theta)$.
- Transform primitives using transformation matrices of SMPL joints and skinning weights.

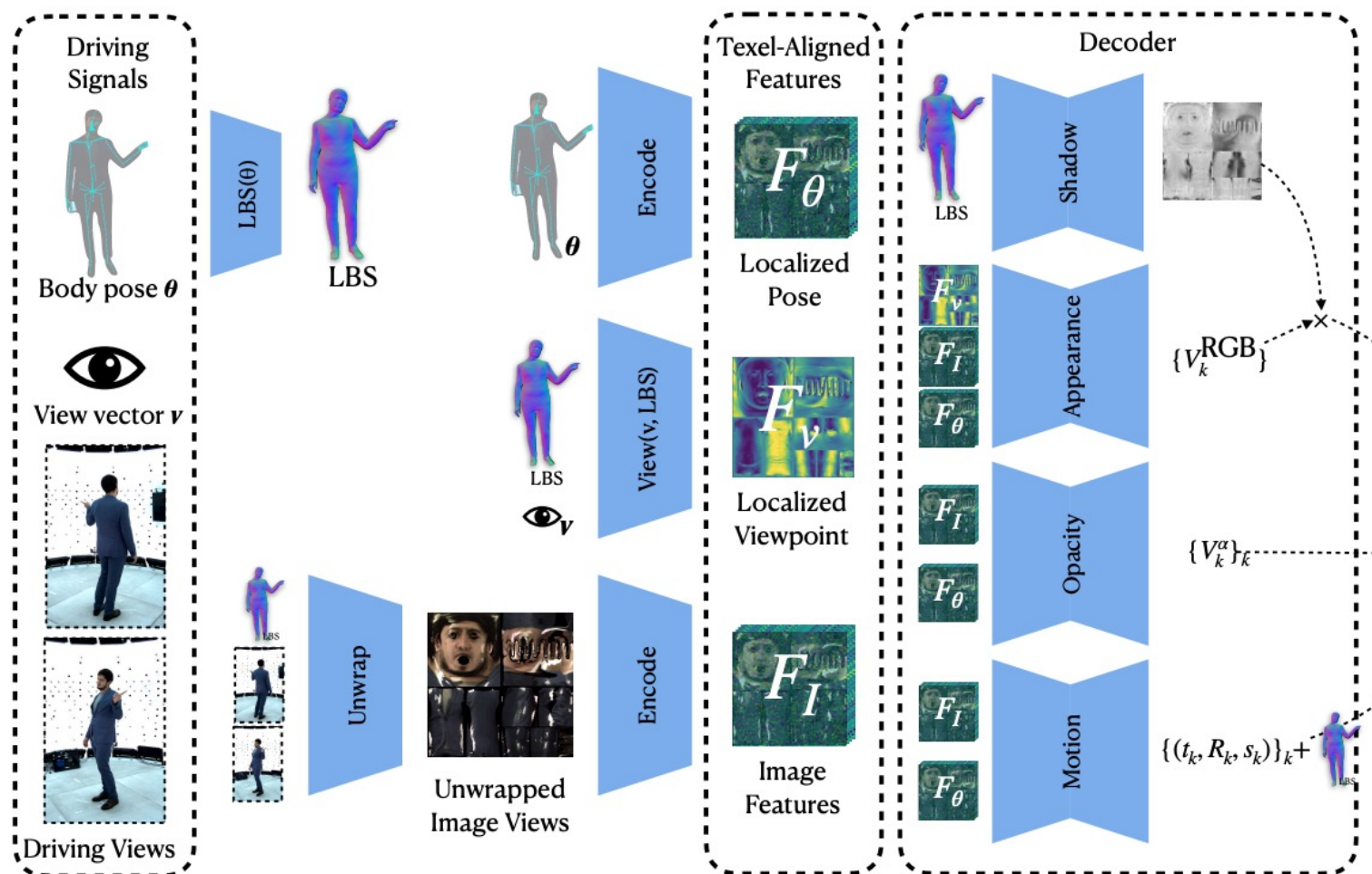
Drivable Volumetric Avatars: Method



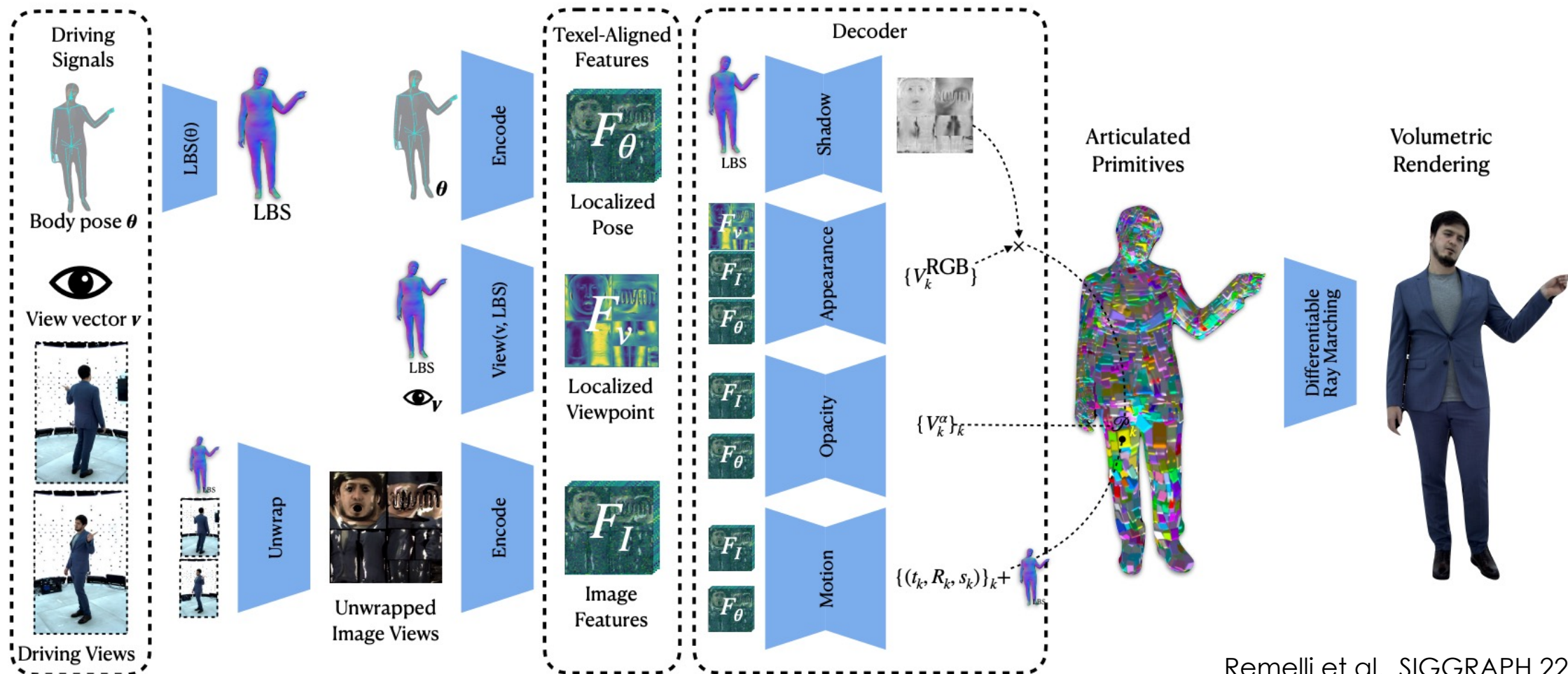
Drivable Volumetric Avatars: Method



Drivable Volumetric Avatars: Method

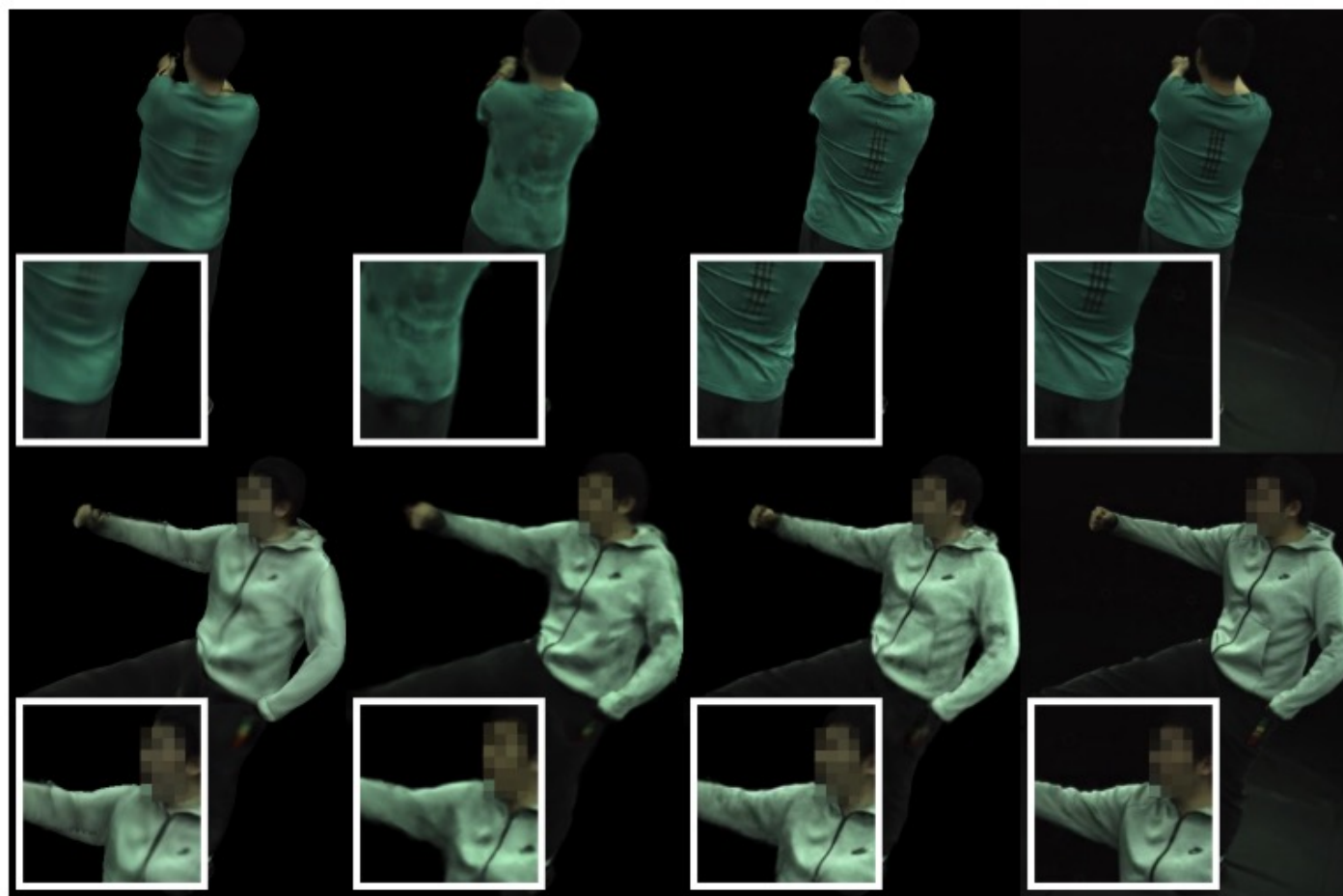


Drivable Volumetric Avatars: Method



Remelli et al. SIGGRAPH 22

Drivable Volumetric Avatars: Results



FBCA

NeuralBody

OURS

Ground Truth

Remelli et al. SIGGRAPH 22

Conclusion:

- For NVS of humans, it helps to introduce human shape and structure prior in the method.
 - Provides controllability
 - Preserves human shape/structure
- Mixture of volumetric primitives helps:
 - Efficient rendering
 - Preserving fine details and model translucent structures

More on Human and NeRFs

- Animatable NeRF, Peng et al., ICCV2021
- H-NeRF, Xu et al., NeurIPS 2021
- NeuMan, Jian et al., ECCV 2022
- DoubleField, Shao et al., CVPR 2022
-
- And many more