

## Motivations & Contributions

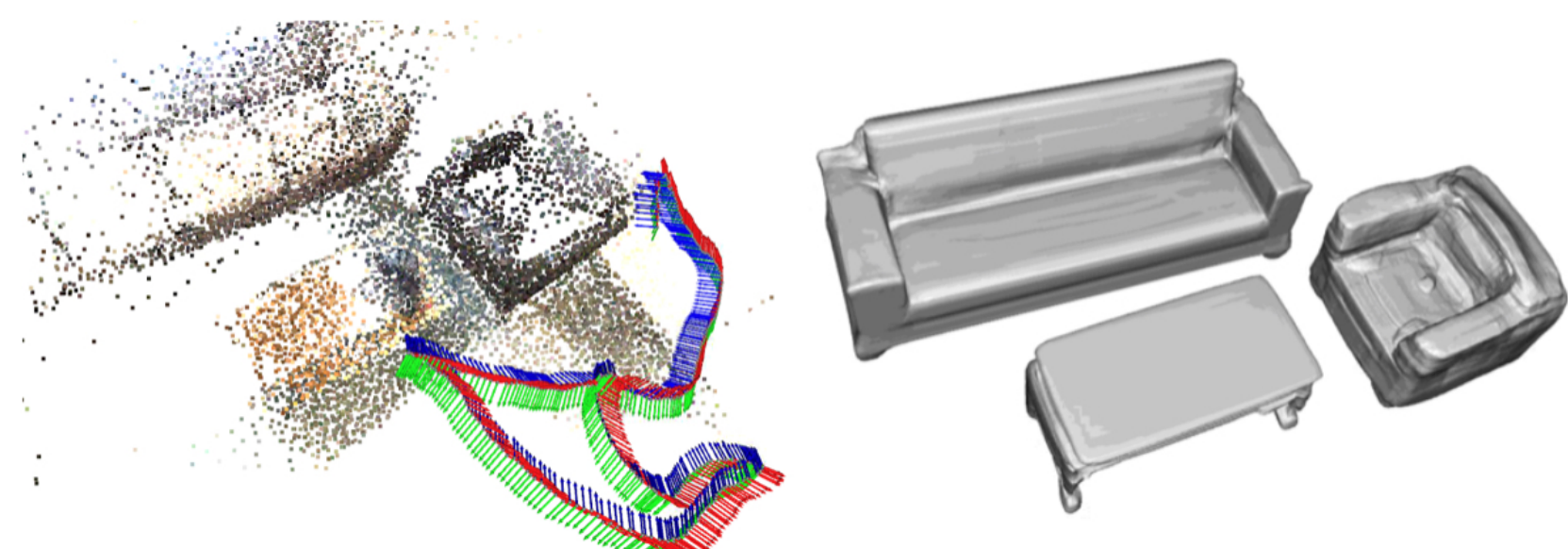
### Motivations:

- Maps offering geometric and semantic information are useful and understandable for humans, allowing specification of symbolic tasks in terms of object entities.
- Striking the right balance between a faithful object reconstruction and a compact object representation remains an open research problem.

### Contributions:

- A bi-level object model consisting of coarse and fine levels, which enables joint optimization of object pose and shape. The coarse-level uses a primitive shape for robust pose and scale initialization, and the fine-level model uses SDF residual directly to allow accurate shape modeling. The two levels are coupled via a shared latent space.
- A cost function to measure the mismatch between the bi-level object model and the segmented RGB-D observations in the world frame.

**Overview:** We propose ELLIPSDF, an expressive yet compact model of object pose and shape, and an associated optimization algorithm to infer an object-level map from multi-view RGB-D camera observations.

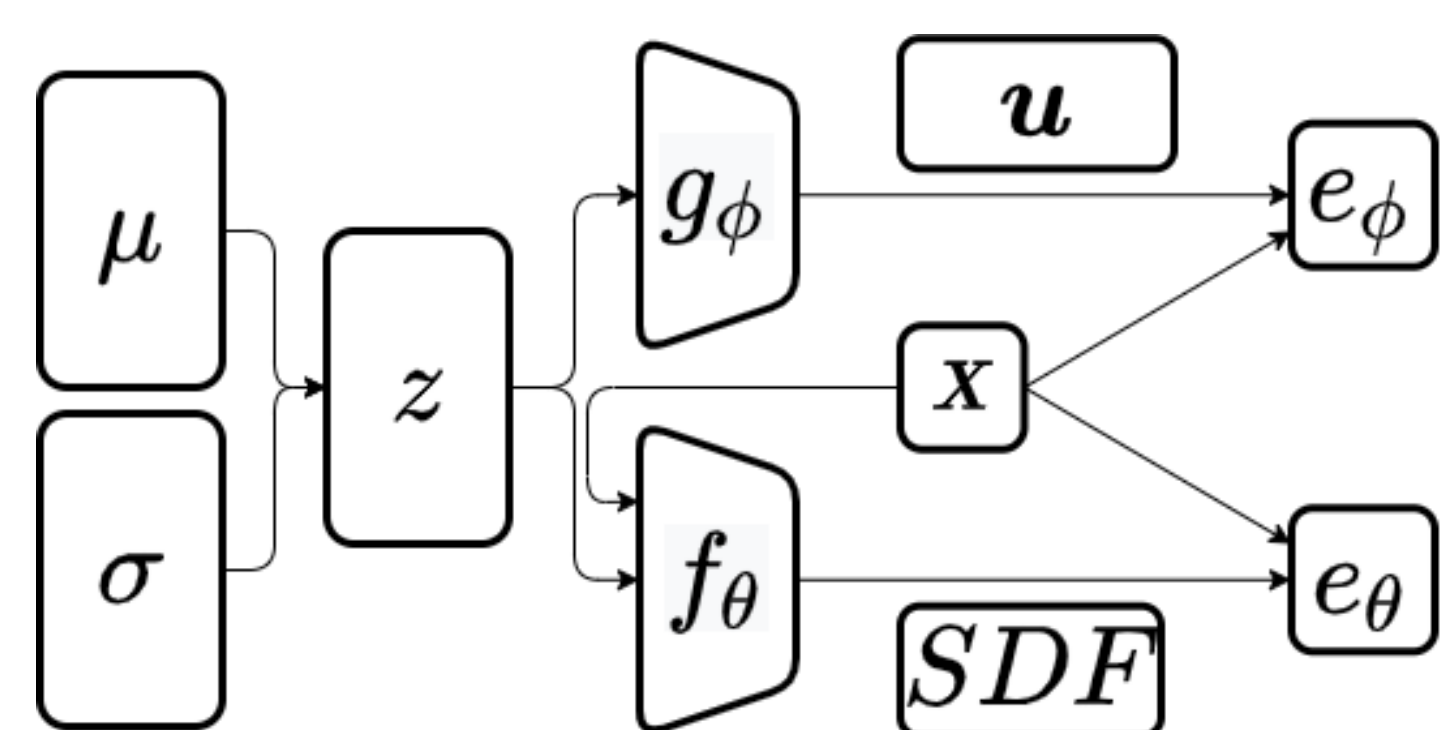


## Object Pose and Shape Optimization

We distinguish between a training phase, where we optimize the parameters  $\mathbf{z}$ ,  $\boldsymbol{\theta}$ ,  $\phi$  of an object class using offline data from instances with known mesh shapes, and a testing phase, where we optimize the pose  $\mathbf{T}$  and shape deformation  $\delta\mathbf{z}$  of a previously unseen instance from the same category using online distance data from an RGB-D camera.

### Training an ELLIPSDF Model:

- Latent shape code shared by coarse shape decoder  $g_\phi$  and fine shape decoder  $f_\theta$ :



- Fine-level shape error function  $e_\theta(\mathbf{x}, d, \mathbf{T}, \delta\mathbf{z}) \triangleq \rho(sf_\theta(\mathbf{PT}\mathbf{x}; \mathbf{z} + \delta\mathbf{z}) - d)$ . Coarse-level shape error  $e_\phi(\mathbf{x}, d, \mathbf{T}, \delta\mathbf{z}) \triangleq \rho(sh(\mathbf{PT}\mathbf{x}, g_\phi(\mathbf{z} + \delta\mathbf{z})) - d)$ .

## Problem Formulation

### Definitions:

- An *object class* is a tuple  $\mathbf{c} \triangleq (\nu, \mathbf{z}, f_\theta, g_\phi)$ , where  $\nu \in \mathbb{N}$  is the class identity, e.g., chair, table, sofa, and  $\mathbf{z} \in \mathbb{R}^d$  is a latent code vector, encoding the average class shape. Shape is represented in a canonical coordinate frame at two levels of granularity: coarse and fine.
- Coarse shape is specified by an ellipsoid  $\mathcal{E}_u$  with semi-axis lengths  $\mathbf{u} = g_\phi(\mathbf{z})$  decoded from the latent code  $\mathbf{z}$  via a function  $g_\phi: \mathbb{R}^d \mapsto \mathbb{R}^3$  with parameters  $\phi$ .
- Fine shape is specified by the signed distance  $f_\theta(\mathbf{x}, \mathbf{z})$  from any  $\mathbf{x} \in \mathbb{R}^3$  to the average shape surface, decoded from the latent code  $\mathbf{z}$  via a function  $f_\theta: \mathbb{R}^3 \times \mathbb{R}^d \mapsto \mathbb{R}$  with parameters  $\theta$ .
- An *object instance* of class  $\mathbf{c}$  is a tuple  $\mathbf{i} \triangleq (\mathbf{T}, \delta\mathbf{z})$ , where  $\mathbf{T} \in \text{SIM}(3)$  specifies the transformation from the global frame to the object instance frame, and  $\delta\mathbf{z} \in \mathbb{R}^d$  is a deformation of the latent code  $\mathbf{z}$ , specifying the average shape of class  $\mathbf{c}$ .

### Error Functions:

- Error function  $e_\phi$  measures the discrepancy between a distance-labelled point  $(\mathbf{x}, d) \in \mathcal{X}_k(\mathbf{p})$  observed close to the instance surface and the coarse shape  $\mathcal{E}_u$  provided by  $\mathbf{u} = g_\phi(\mathbf{z})$ . Error function  $e_\theta$  is used for the difference between  $(\mathbf{x}, d)$  and the SDF value  $f_\theta(\mathbf{x}, \mathbf{z})$  predicted by the fine shape model.
- Overall error function is defined as:

$$\alpha e_r(\delta\mathbf{z}) + \sum_k \sum_{\mathbf{p}} \sum_{(\mathbf{x}, d)} \beta e_\theta(\mathbf{x}, d, \mathbf{T}, \delta\mathbf{z}) + \gamma e_\phi(\mathbf{x}, d, \mathbf{T}, \delta\mathbf{z}).$$

### Joint Pose and Shape Optimization:

- The Jacobian of  $e_\theta$  wrt transformation perturbation is:

$$\frac{\partial e_\theta}{\partial \xi} = \frac{\partial \rho(r)}{\partial r} \left( \hat{s}[0_6, 1] f_\theta(\mathbf{x}, \delta\mathbf{z}) + \hat{s} \nabla_{\mathbf{x}} f_\theta(\mathbf{x}, \delta\mathbf{z})^\top \mathbf{P} [\hat{\mathbf{T}} \mathbf{x}]^\odot \right)$$

$$\frac{\partial e_\theta}{\partial \delta\mathbf{z}} = \frac{\partial \rho(r)}{\partial r} \hat{s} \nabla_{\mathbf{z}} f_\theta(\mathbf{x}, \delta\mathbf{z}).$$

- Given initial transformation and deformation, solve joint pose and shape optimization via gradient descent:

$$\mathbf{T}^{i+1} \triangleq \exp \left( -\eta_1 \frac{\partial e(\mathbf{T}, \delta\mathbf{z}, \boldsymbol{\theta}^*, \boldsymbol{\phi}^*; \{\mathcal{X}_k(\mathbf{p})\})}{\partial \mathbf{x} \mathbf{i}} \right) \mathbf{T}^i,$$

$$\delta\mathbf{z}^{i+1} \triangleq \delta\mathbf{z}^i - \eta_2 \left( \frac{\partial e(\mathbf{T}, \delta\mathbf{z}, \boldsymbol{\theta}^*, \boldsymbol{\phi}^*; \{\mathcal{X}_k(\mathbf{p})\})}{\partial \delta\mathbf{z}} \right).$$

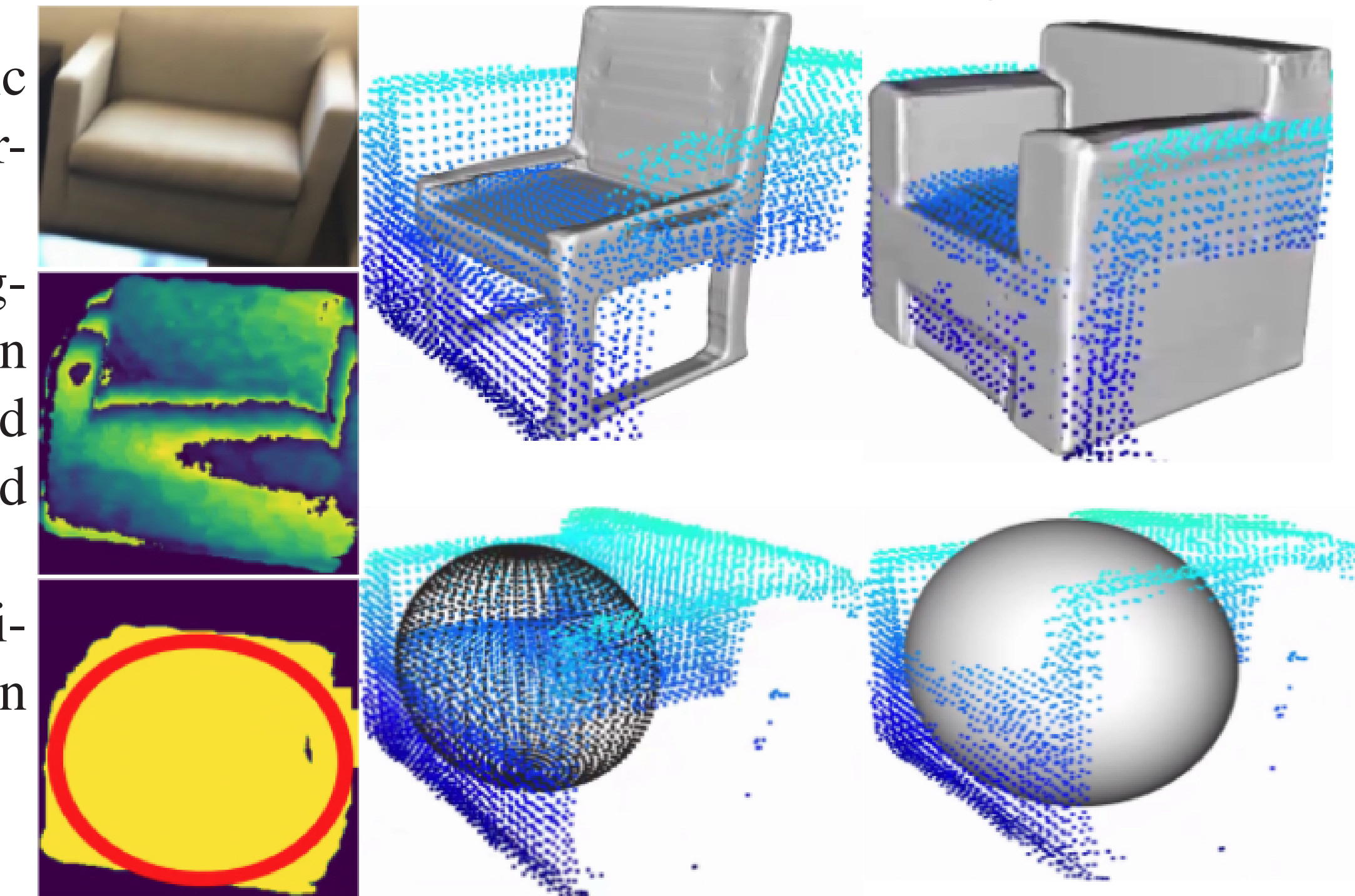
## Experiments & Results

We evaluate ELLIPSDF on the ScanNet dataset, which provides 3D scans captured by a RGB-D sensor of indoor scenes with chairs, tables, displays etc. We segment out the objects from the scene-level mesh using provided instance labels and sample points from the object meshes to generate the point observations.

### Visualizations of Intermediate Results:

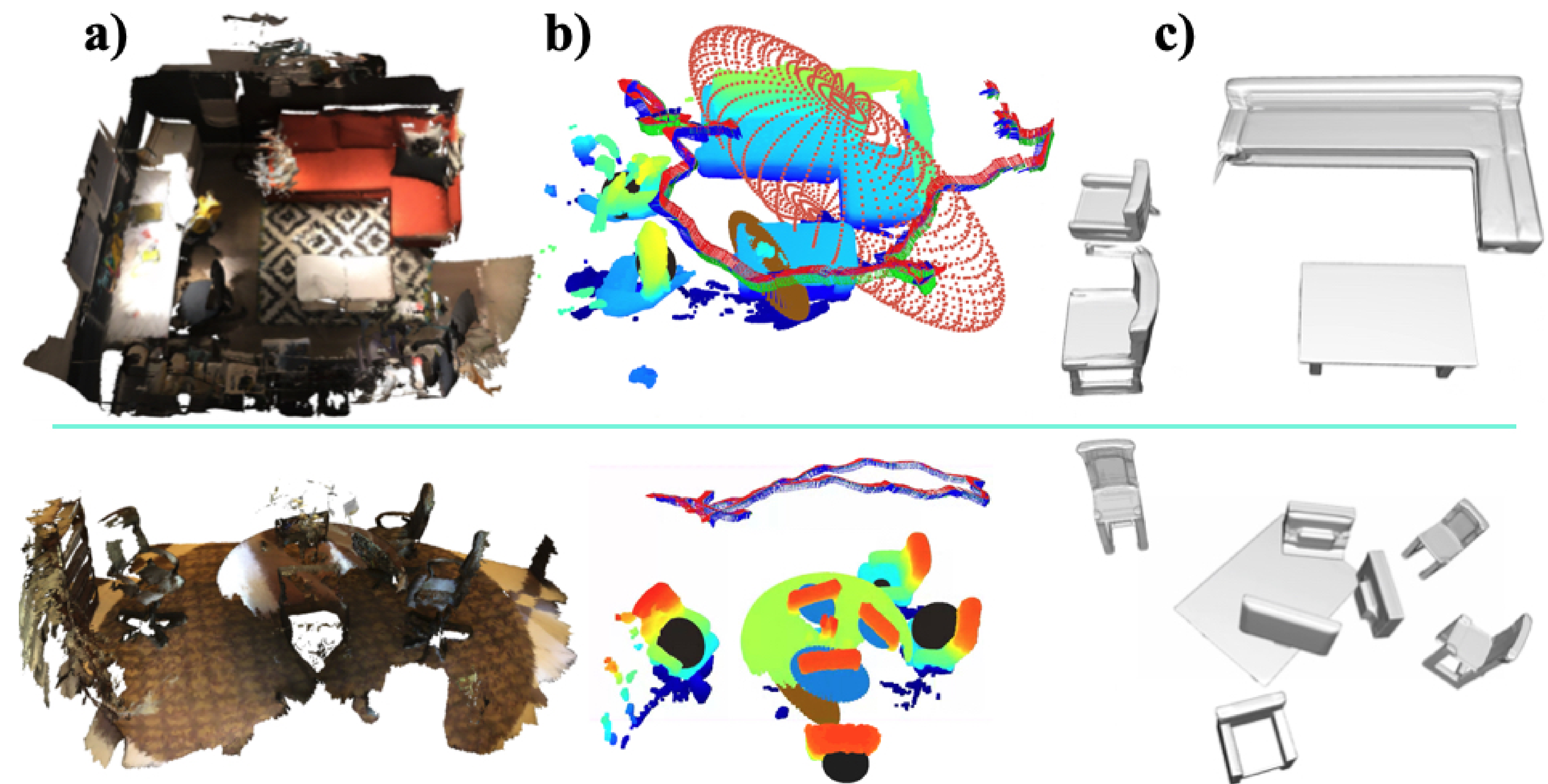
- The ELLIPSDF decoder model is trained on synthetic CAD models from ShapeNet. Each model's scale is normalized to be inside a unit sphere.
- First column: RGB image, depth image, instance segmentation (yellow), fitted ellipse (red) for a chair in ScanNet scene 0461. Second column: mean shape and ellipsoid with initialized pose. Third column: optimized fine-level and coarse-level shapes with optimized pose.
- Optimization step improves the scale and shape estimates notably, e.g. by transforming the four-leg mean shape into an armchair.

Intermediate ELLIPSDF stages.



### Qualitative Results on a larger scale:

Column a): Ground-truth scene in ScanNet Sequences. Column b): The ellipsoids (black for chair, red for sofa, blue for monitor, brown for table) are the initialized objects. Column c): Reconstructed meshes using ELLIPSDF.



### Quantitative results for pose estimation on ScanNet:

Scan2CAD	Vid2CAD	ELLIPSDF (init)	ELLIPSDF (opt)
31.7	38.3	31.5	<b>39.6</b>

### Quantitative results for shape evaluation on ScanNet:

Method	cabinet	chair	display	table	avg.
# instances	132	820	209	146	327
ELLIPSDF (fine)	88.4	88.3	90.6	76.2	85.9
ELLIPSDF (coarse+fine)	<b>91.0</b>	<b>90.6</b>	<b>96.9</b>	<b>77.3</b>	<b>89.0</b>

### Comparison of 3D detection results on ScanNet:

mAP @ IoU=0.5	Chair	Table	Display
FroDO	0.32	0.06	0.04
MOLTR	0.39	0.06	0.10
ELLIPSDF (fine)	0.42	0.26	0.25
ELLIPSDF (coarse+fine)	<b>0.43</b>	<b>0.27</b>	<b>0.31</b>