



## Take away

- The first **single-view self-supervised method** for depth estimation.
- Light and camera co-located + dark environment (e.g.: **endoscopy**)
- We outperform multi-view self-supervision and **match supervision with ground truth**

## Motivation

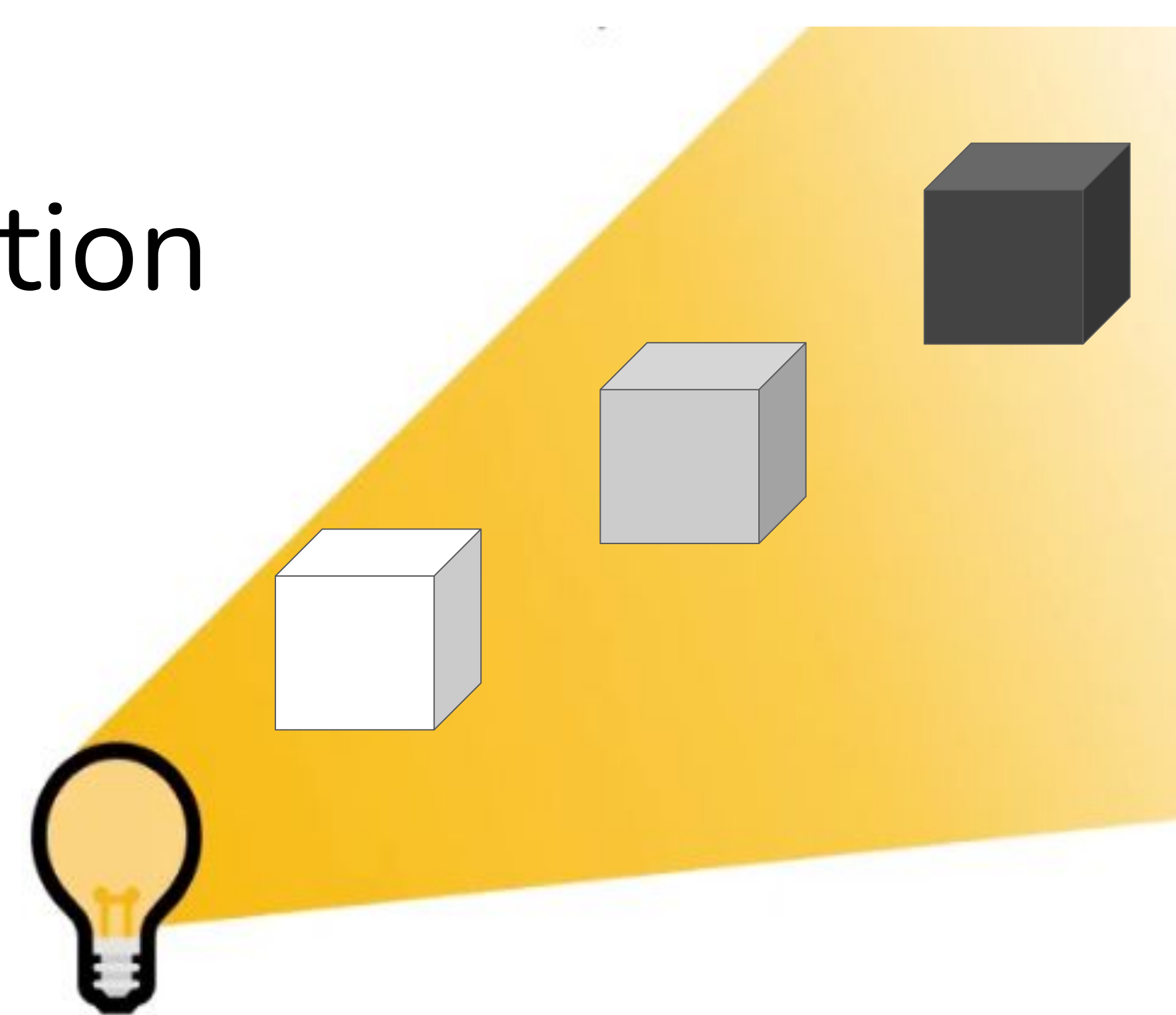
Single-View Depth in Endoscopy

- Monocular camera
- Varying illumination
- Untextured surface
- Deformable environment

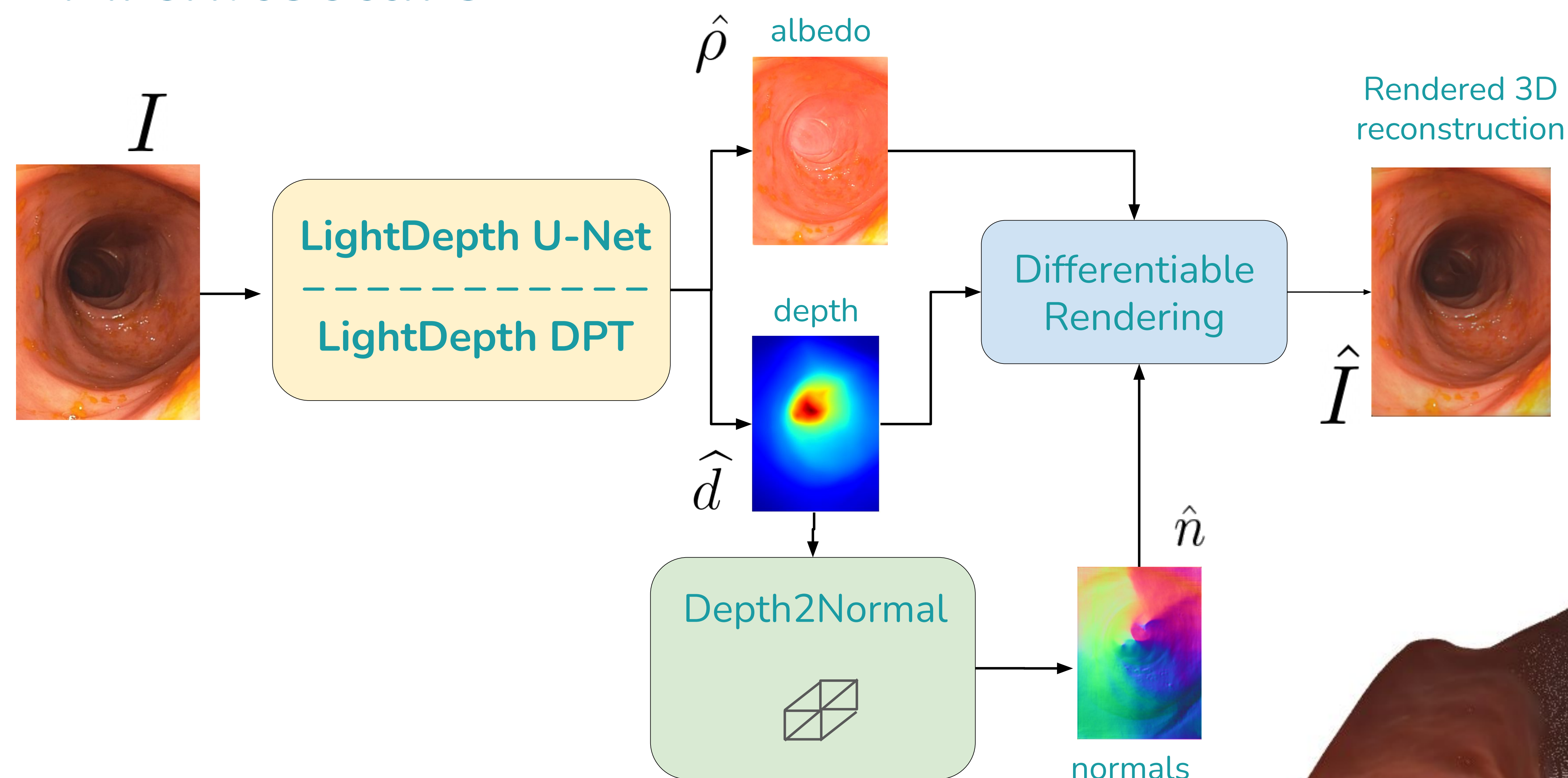


## Illumination decline as single-view self-supervision

- The further, the darker
- No need for depth GT
- No need for camera motion estimation
- Single architecture for training and test-time refinement (TTR)



## Architecture



## Differentiable rendering

$$\mathcal{I}(d_i, \rho_i, g) = \left( \frac{\sigma_0}{\|d_i \mathbf{r}_i - \mathbf{x}_l\|^2} R(\psi_i) \cos(\theta_i) \rho_i g \right)^{1/\gamma}$$

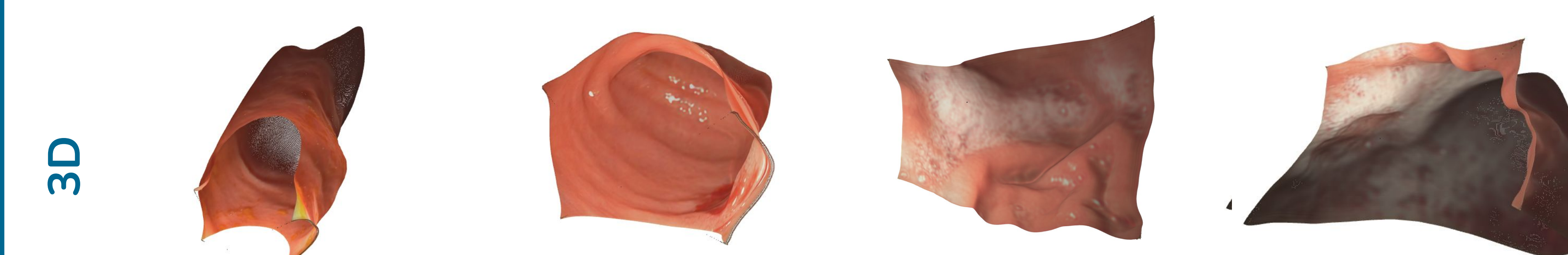
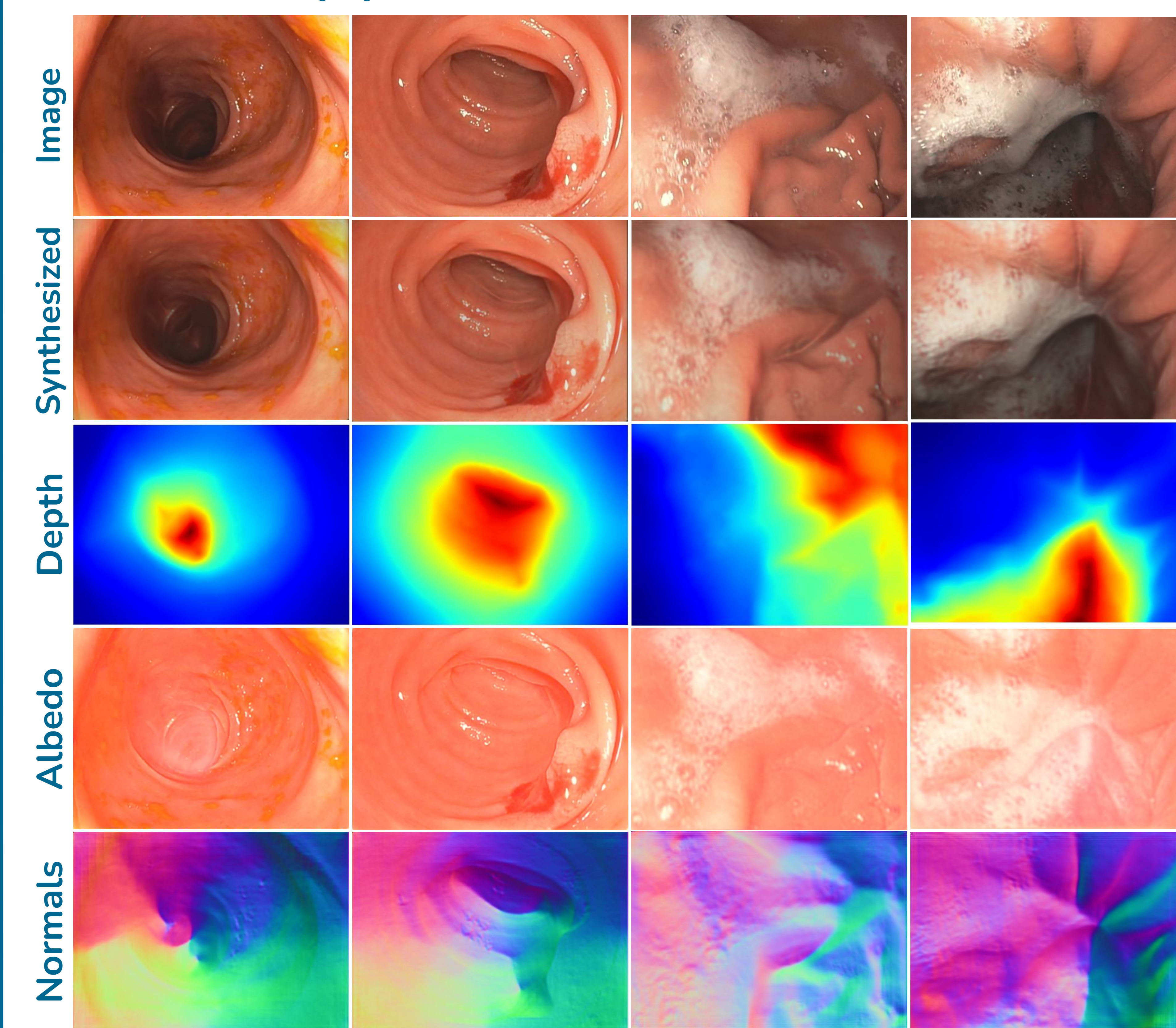
## Losses

- Single-view photometric consistency  $\mathcal{L}_p = \sum_{i \in \Omega} (I_i - \hat{I}_i)^2$
- Specularity consistency  $\mathcal{L}_{sp} = \sum_{i \in \Omega} (\cos \alpha_i - 1)^2$
- Smoothness  $\mathcal{L}_s = |\partial_x \hat{d}| e^{-|\partial_x I|} + |\partial_y \hat{d}| e^{-|\partial_y I|}$

## Experiments

- Metrics match supervised methods
- Significant better performance than multi-view self-supervision

## EndoMapper dataset



Architecture	Backbone	Supervision	MAE ↓	MedAE ↓	RMSE ↓
U-Net	ResNet18	Depth GT	4.15	3.29	5.52
DPT-Hybrid [48]	ResNet50	Depth GT	<b>3.22</b>	2.77	<b>4.10</b>
Monodepth2 [20]	ResNet50	Multi-View	14.27	9.59	18.64
CADepth [64]	ResNet18	Multi-View	52.35	17.04	87.43
XDCycleGAN [42]	ResNet	Cycle	17.16	11.91	22.43
LightDepth U-Net	ResNet18	Light	4.37	2.92	6.31
LightDepth DPT	ResNet50	Light	3.94	2.67	5.60
LightDepth U-Net	ResNet18	Light (TTR)	3.72	2.59	5.43
LightDepth DPT	ResNet50	Light (TTR)	<b>3.70</b>	<b>2.58</b>	<b>5.27</b>