

Self-Supervised Monocular Scene Decomposition and Depth Estimation

Sadra Safadoust, Fatma Güney

KUIS AI Center, Koç University



Self-Supervised Monocular Depth

Current monocular depth estimation methods

- either assume a static scene and fail in foreground regions with independently moving objects
- or require a separate segmentation step to identify the dynamic objects in the foreground.

MonoDepthSeg

We introduce **MonoDepthSeg** to jointly estimate depth and segment moving objects from monocular video without using any ground-truth labels.



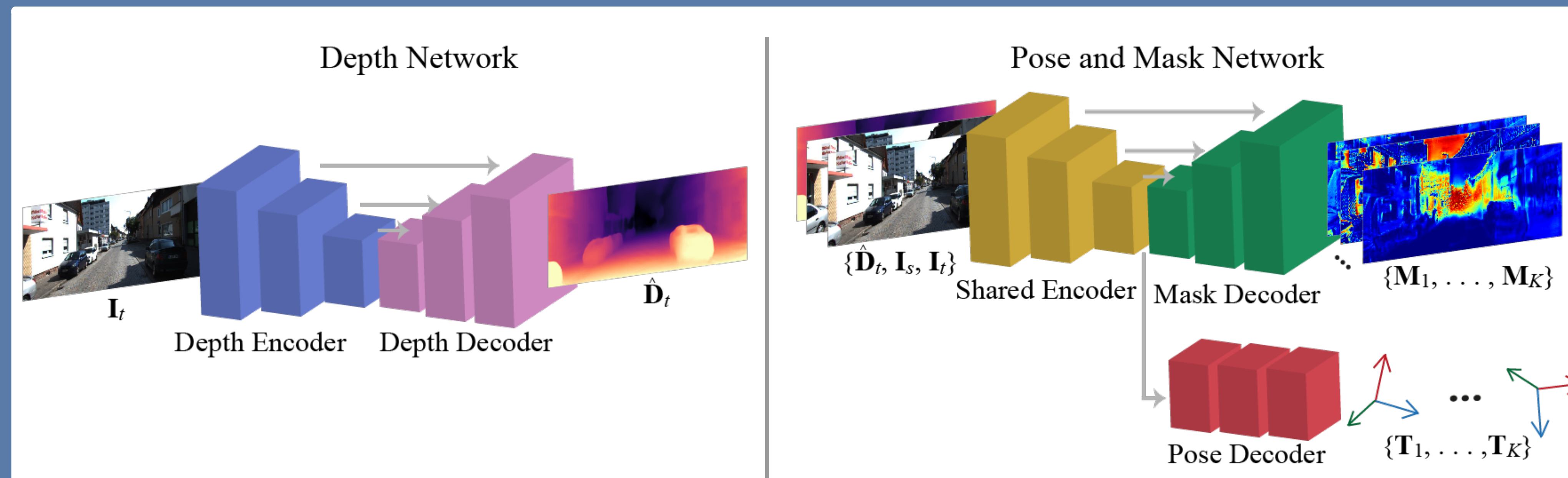
We decompose the scene into a fixed number of components where each component corresponds to a region on the image with its own transformation matrix representing its motion. This improves the results in regions with moving objects (bottom-right) compared to the current approaches [1] (bottom-left), while simultaneously recovering a decomposition of the scene, mostly corresponding to moving regions (top-right).

Our framework consists of

- depth network to estimate per-pixel depth values
- pose and mask network to divide the image into components and estimate a separate pose for each component.

The two networks are trained jointly and end-to-end.

Model Architecture



- Given a single target image I_t , the depth network outputs the depth estimate \hat{D}_t .
- The shared encoder maps the two consecutive input frames I_s and I_t and the depth estimate \hat{D}_t to a common representation.
- The mask decoder produces the same resolution K masks $\{M_1, \dots, M_K\}$.
- The pose decoder maps the same encoded representation into rigid transformations $\{T_1, \dots, T_K\}$ corresponding to the masks.

Methodology

- We represent the motion of each component using a 3D rigid transformation $T = [R, t] \in SE(3)$
- We encourage masks to be layered according to a pre-defined depth order d_i . This helps to account for occlusions:

$$M_i(p) = \frac{e^{d_i M'_i(p)}}{\sum_{j=1}^K e^{d_j M'_j(p)}}$$
- For every pixel p on the target image I_t , we compute the corresponding 3D point x using its depth value \hat{D}_t and the intrinsic camera matrix K :

$$x = \hat{D}_t(p) K^{-1} p$$
- We transform the 3D point x using the masks and rigid transformations to obtain x' :

$$x' = \sum_{i=1}^K M_i(p) T_i x = \sum_{i=1}^K M_i(p) (R_i x + t_i)$$
- We project the transformed point x' to find the corresponding point p' on the source image I_s :

$$p' = K x'$$
- We reconstruct the target image I_t by sampling pixels from the source image I_s and obtain the warped image \hat{I}_s , such that $\hat{I}_s(p) = I_s(p')$.
- We use an edge-aware smoothness loss \mathcal{L}_{smooth} over the mean-normalized inverse depth values and define the photometric loss as follows:

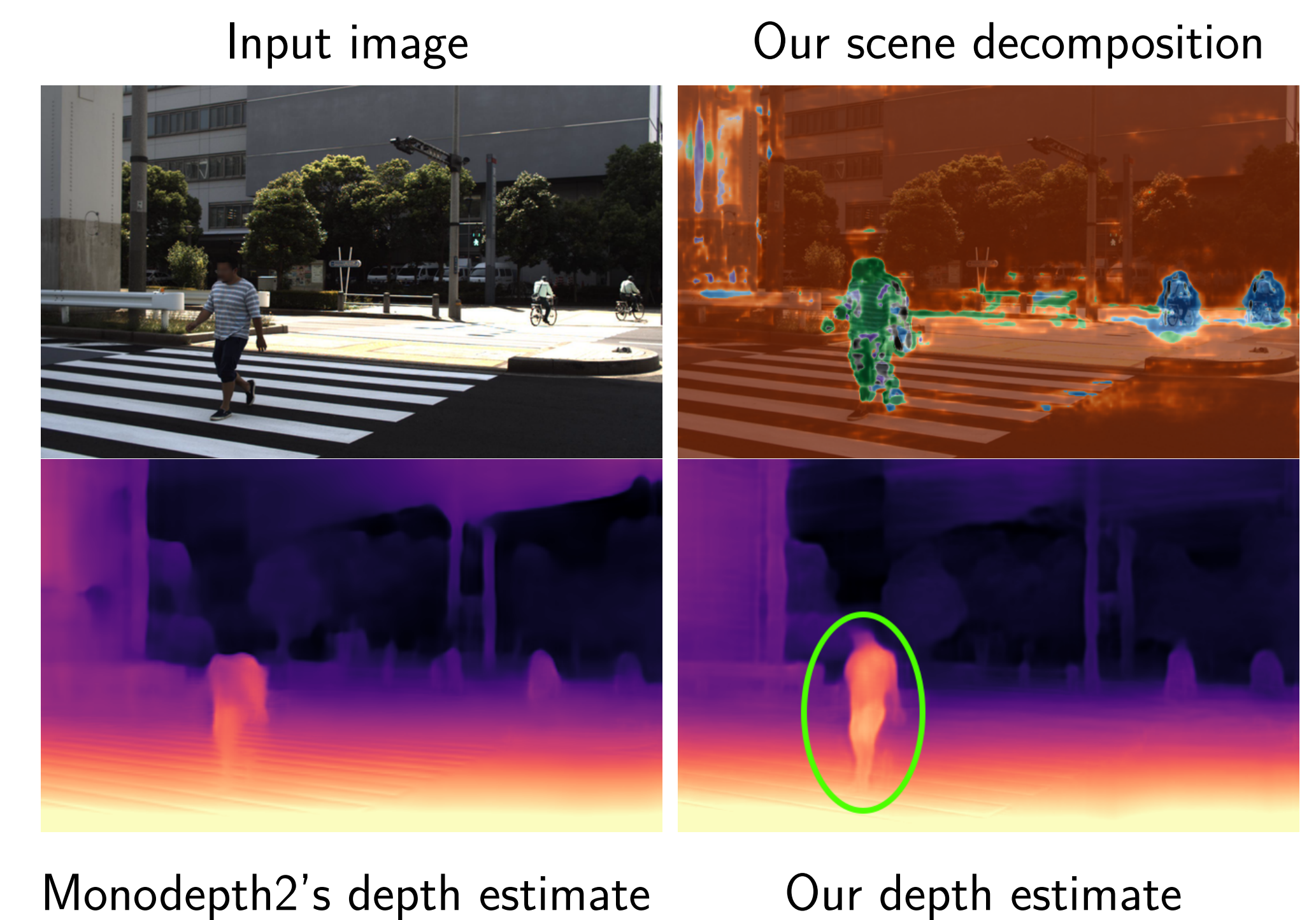
$$\mathcal{L}_{photo}(p) = \min_s \left[(1 - \alpha) |I_t(p) - \hat{I}_s(p)| + \frac{\alpha}{2} \left(1 - SSIM(I_t, \hat{I}_s)(p) \right) \right]$$
- Our final loss is defined as:

$$\mathcal{L} = \frac{1}{N} \sum_p \mathcal{L}_{photo}(p) + \lambda \mathcal{L}_{smooth}(p)$$

Quantitative Results

	Method	Abs Rel	RMSE	Car	Person
DDAD	PackNet [2]	0.23	17.92	0.38	0.20
	Monodepth2 [1]	0.22	17.63	0.25	0.21
	Ours	0.19	16.61	0.24	0.17
	Method	Abs Rel	RMSE		
CityScapes		Moving	All	Moving	All
	Monodepth2 [1]	0.158	0.170	8.043	8.155
	Ours	0.143	0.142	7.649	7.361
KITTI	Monodepth2 [1]	0.143	0.110	5.949	4.642
	Ours	0.138	0.110	5.796	4.700

Qualitative Results



References

- [1] C. Godard, O. Mac Aodha, M. Firman, and G. J. Brostow, "Digging into self-supervised monocular depth estimation," 2019.
- [2] V. Guizilini, R. Ambrus, S. Pillai, A. Raventos, and A. Gaidon, "3D packing for self-supervised monocular depth estimation," 2020.

Contact Information

Web: <https://kuis-ai.github.io/monodepthseg>
 Email: {ssafadoust20, fgüney}@ku.edu.tr