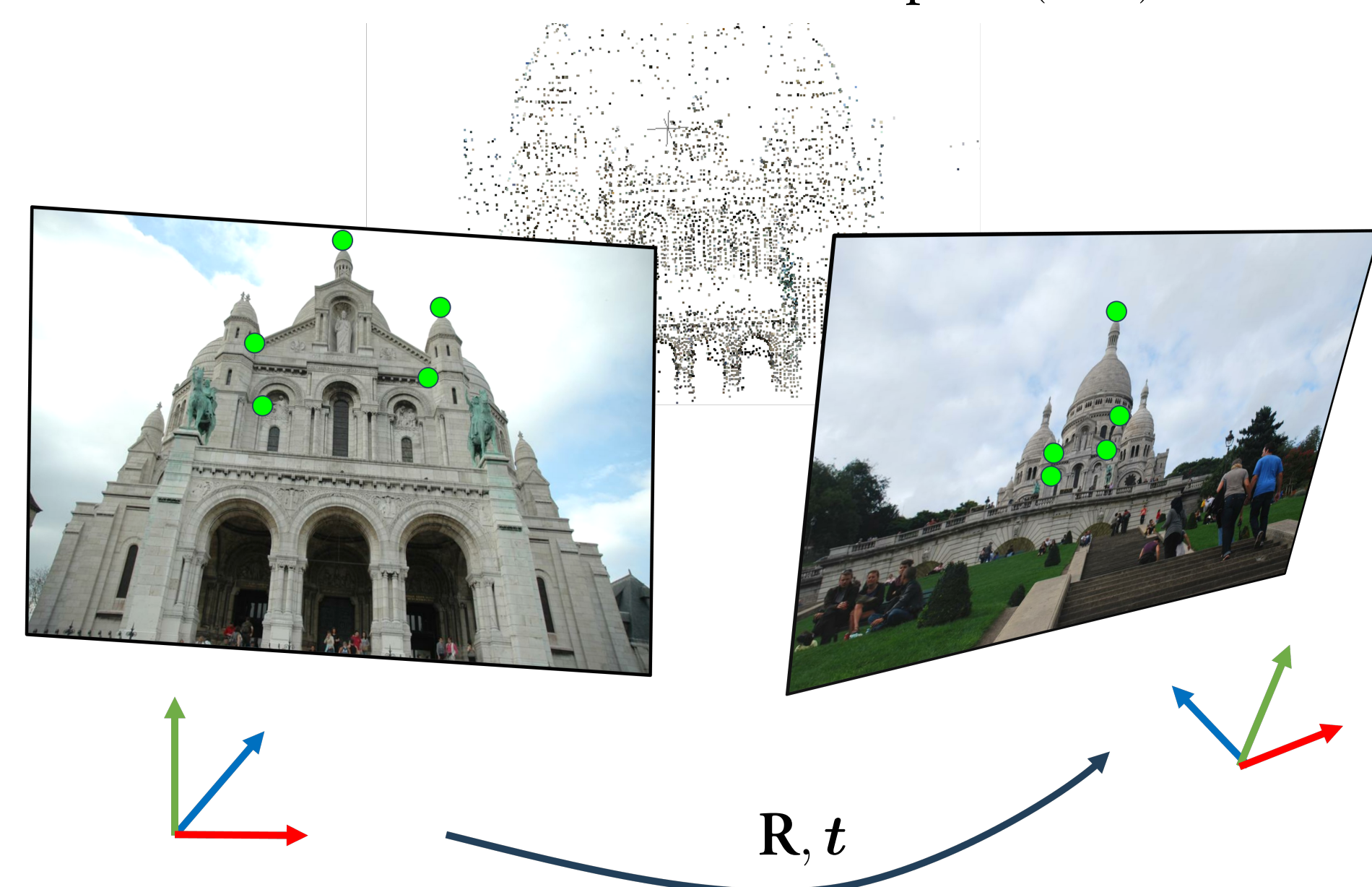


Summary

Goal: Improve estimation of the relative camera pose (\mathbf{R}, \mathbf{t}) between two images.



- Given a sparse set of keypoint correspondences, the relative camera pose can be estimated using RANSAC.
- For each point-correspondence, in addition to the positions (x, y) , we use the *relative depth*, i.e. relative distance to the same scene point in the two images.
- Using this extra constraint we can generate pose candidates for RANSAC using fewer point correspondences, compared to purely coordinate-based solvers.

Contributions

- A novel 3-point minimal solver for relative pose, using relative depths.
- We show that the relative depth can either be estimated from SIFT scales, or predicted using a simple neural network.
- Through experiments, we demonstrate that the smaller sample size leads to a significantly reduced runtime in settings with high outlier ratios, compared to purely point-based solvers.

Relative Pose Estimation

The projections \mathbf{x}, \mathbf{x}' of a 3D-point \mathbf{X} are described by the camera equations

$$\begin{cases} \lambda \mathbf{x} = \mathbf{X} \\ \lambda' \mathbf{x}' = \mathbf{R}\mathbf{X} + \mathbf{t} \end{cases} \Rightarrow \lambda' \mathbf{x}' = \lambda \mathbf{R}\mathbf{x} + \mathbf{t}, \quad (1)$$

where λ and λ' are the depths of point \mathbf{X} .

- Classical minimal solver **requires 5 points** to estimate relative pose.
- In RANSAC, number of iterations grows exponentially with sample size.

Relative Depth in Relative Pose Estimation

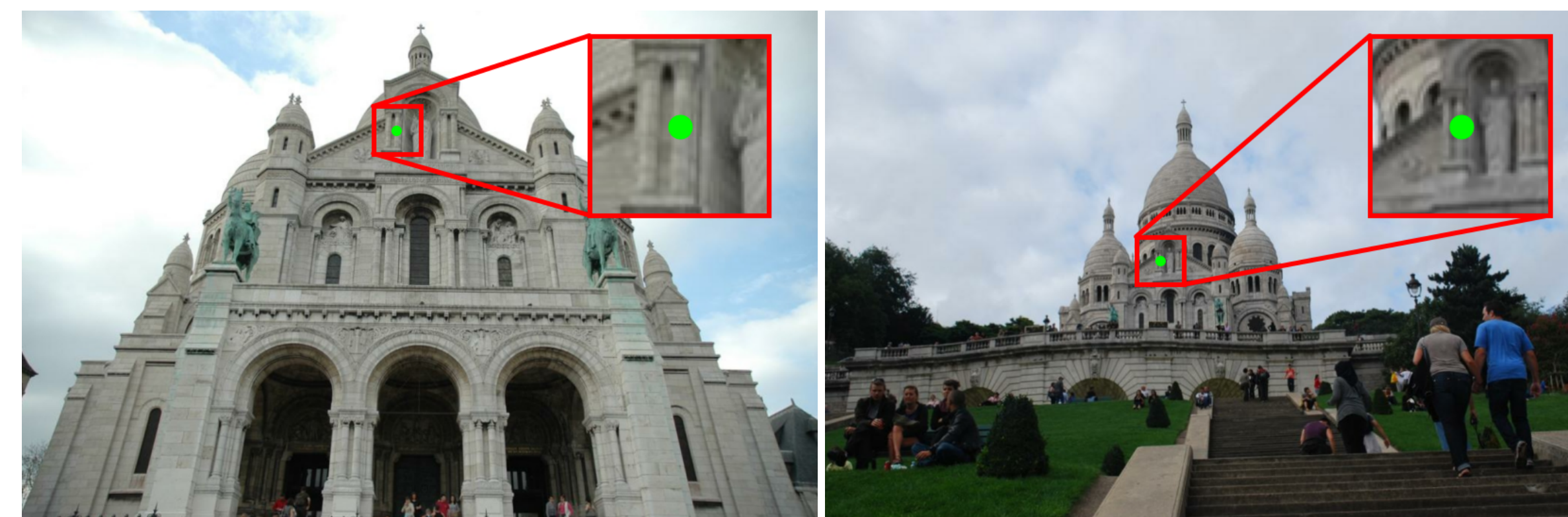
Idea: Leverage relative depth constraints, observed from scale changes.

- If we introduce relative depth $\sigma := \lambda'/\lambda$, we can rewrite (1) as

$$\lambda(\sigma \mathbf{x}' - \mathbf{R}\mathbf{x}) = \mathbf{t}. \quad (2)$$

- Relative depth inversely proportional to the relative scale in the images

$$\sigma := \frac{\lambda'}{\lambda} = \frac{f' s}{f s'}. \quad (3)$$



- Keypoint detection scale (e.g. from SIFT) can be used directly in (3).
- From (2), we introduce a novel minimal 3-point solver

$$\begin{aligned} \sigma_1 \lambda_1 \mathbf{x}'_1 &= \lambda_1 \mathbf{R}\mathbf{x}_1 + \mathbf{t}, \\ \sigma_2 \lambda_2 \mathbf{x}'_2 &= \lambda_2 \mathbf{R}\mathbf{x}_2 + \mathbf{t}, \\ \lambda_3 \mathbf{x}'_3 &= \lambda_3 \mathbf{R}\mathbf{x}_3 + \mathbf{t}. \end{aligned}$$

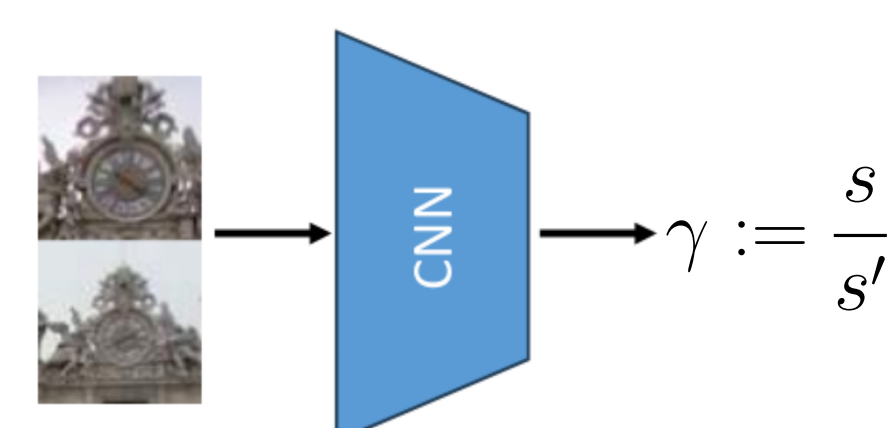
- Forming the differences and taking the norm eliminates \mathbf{R} and yields

$$\begin{aligned} \|\sigma_1 \mathbf{x}'_1 - \sigma_2 \lambda_2 \mathbf{x}'_2\|^2 &= \|\mathbf{x}_1 - \lambda_2 \mathbf{x}_2\|^2, \\ \|\sigma_1 \mathbf{x}'_1 - \lambda_3 \mathbf{x}'_3\|^2 &= \|\mathbf{x}_1 - \lambda_3 \mathbf{x}_3\|^2, \\ \|\sigma_2 \lambda_2 \mathbf{x}'_2 - \lambda_3 \mathbf{x}'_3\|^2 &= \|\lambda_2 \mathbf{x}_2 - \lambda_3 \mathbf{x}_3\|^2. \end{aligned}$$

- We can find the three unknowns (red) by solving two quadratics.
- In the paper we also show extension to known vertical direction (2-points).

RelScaleNet

To improve scale estimate, we introduce a simple neural network that directly regresses relative scale from pairs of image patches.

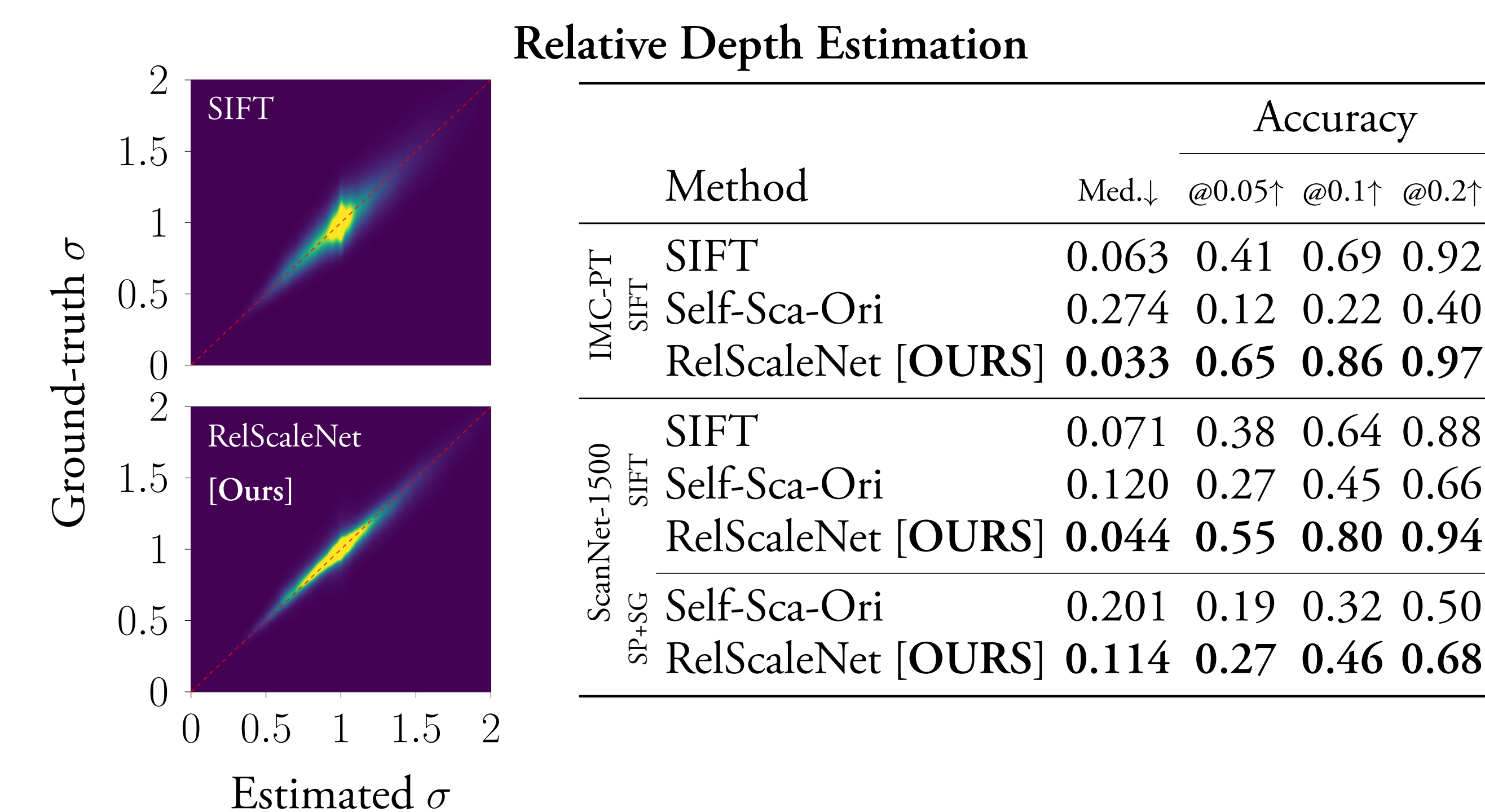


Input: A pair of image patches from neighborhood of corresponding points.

Output: Estimate of relative scale $\gamma = s/s'$.

- We train on MegaDepth and supervise with MSE-loss w.r.t. ground-truth γ .
- Ground-truth relative depth calculated from Structure-from-Motion model.

Evaluation



Relative Pose Estimation with LO/GC-RANSAC on IMC-PT

RSC Method	All pairs		Hardest 5%		
	AUC@5°	RT(ms)	AUC@5°	RT(ms)	
LO-RANSAC	5 pt. (Nistér)	56.89	15.7	12.13	42.2
	3 pt. + SIFT (Barath & Kukulova)	30.77	7.0	1.23	21.3
	3 pt. + SIFT [OURS]	54.30	<u>13.4</u>	8.72	2.8
	3 pt. + RelScaleNet [OURS]	<u>54.63</u>	15.0	<u>9.47</u>	2.8
GC-RANSAC	5 pt. (Nistér)	56.22	25.4	9.76	16.5
	3 pt. + SIFT (Barath & Kukulova)	50.55	11.1	2.16	6.0
	3 pt. + SIFT [OURS]	52.73	<u>16.2</u>	5.24	4.7
	3 pt. + RelScaleNet [OURS]	<u>53.11</u>	16.8	<u>5.65</u>	4.7

Relative Pose Estimation with LO-RANSAC on ScanNet-1500

KP. Method	AUC@5°	AUC@10°	AUC@20°	Runtime (ms)	
SIFT	5 pt. (Nistér)	11.06	21.99	33.32	8.2
	3 pt. + SIFT (Barath & Kukulova)	4.94	10.33	17.16	<u>3.7</u>
	3 pt. + SIFT [OURS]	9.90	20.59	31.96	2.9
	3 pt. + RelScaleNet [OURS]	<u>10.43</u>	<u>21.21</u>	<u>32.43</u>	2.9
SP+SG	5 pt. (Nistér)	<u>17.55</u>	<u>34.21</u>	<u>51.50</u>	<u>59.4</u>
	3 pt. + RelScaleNet [OURS]	18.39	35.46	52.24	10.4

Conclusion

Our novel 3-point solver has similar accuracy to the 5-point solver, while being significantly faster in high outlier settings.