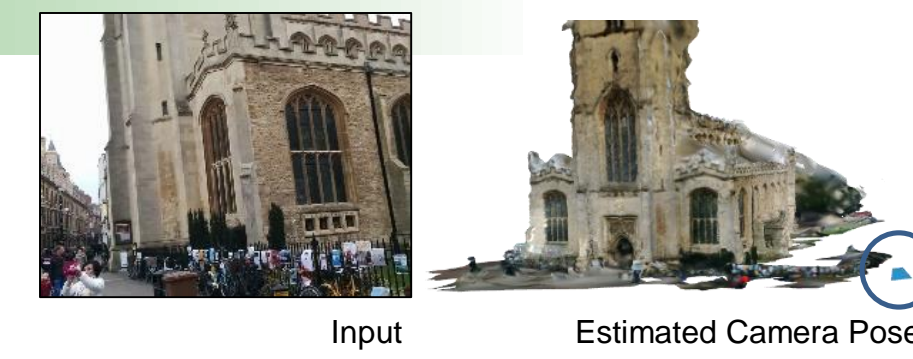




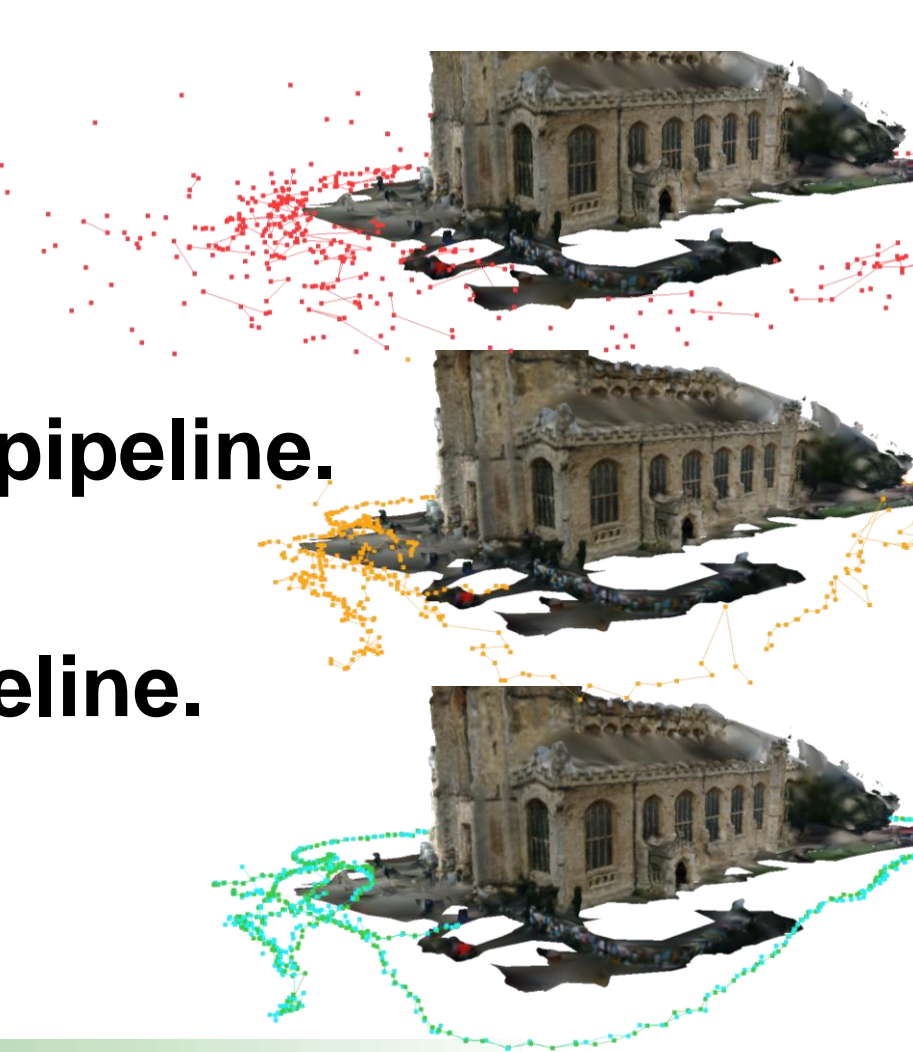
Problem Statement

Estimate the **6D camera pose** (position + orientation) relative to a known scene from a **single RGB image**.



We show that learning less is more. See right:

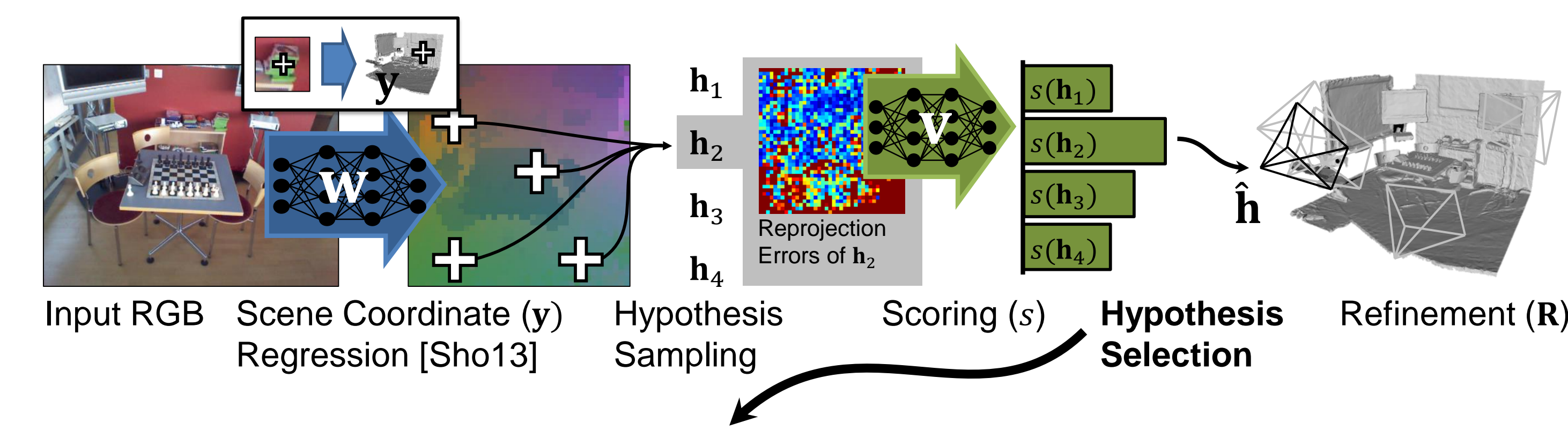
- **Red:** Learn everything. CNN predicts pose directly.
- **Orange:** Learn two components of a geometric pipeline. Our previous work [Bra17].
- **Cyan:** Learn one component of a geometric pipeline. This work.
- **Green:** Ground truth camera path.



Contributions

- **Fully differentiable, robust pose optimization** without learnable parameters **on top of learned scene coordinate regression**
- **Learning** scene coordinate regression **without a 3D scene model or depth maps**
- **Stable end-to-end training** due to new approximation of refinement gradients, and controlling the entropy of pose hypotheses
- We **exceed state-of-the-art** on camera localization **on three datasets** (indoor and outdoor)

Previous Work: Differentiable RANSAC (DSAC) [Bra17]



Probabilistic Pose Selection

$$\hat{\mathbf{h}}^w = \mathbf{h}_j^w, \text{ where } j \sim P(j|\mathbf{w})$$

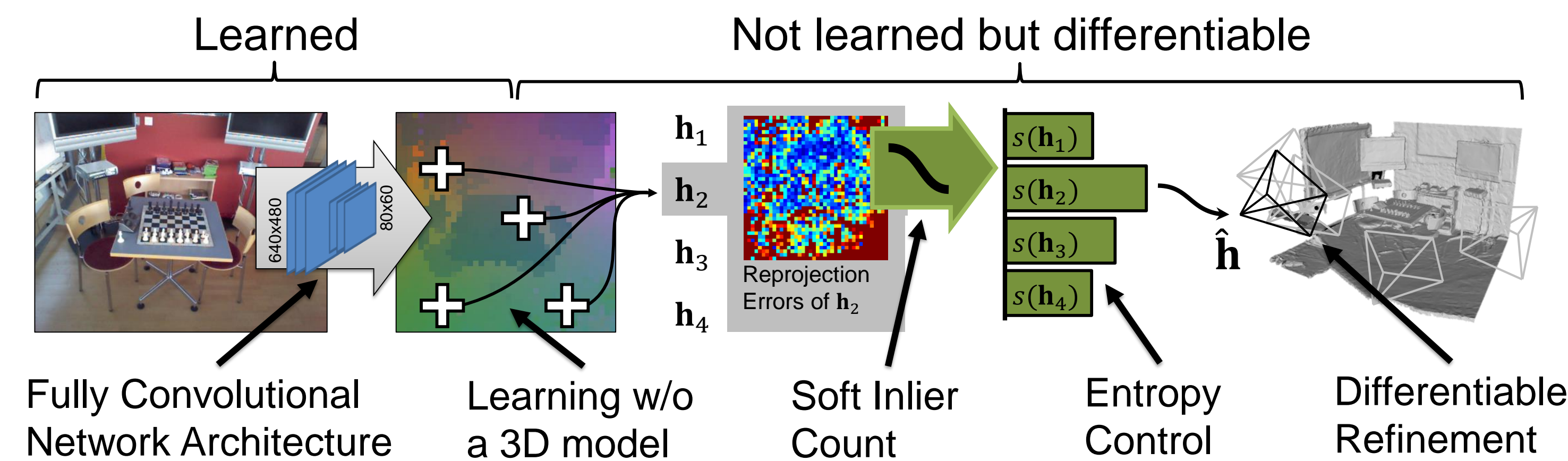
$$P(j|\mathbf{w}) = \frac{\exp(s(\mathbf{h}_j^w))}{\sum_k \exp(s(\mathbf{h}_k^w))}$$

DSAC Learning Objective

$$\frac{\partial}{\partial \mathbf{w}} \mathbb{E}_{j \sim P(j|\mathbf{w}, \mathbf{v})} [\ell(\mathbf{R}(\mathbf{h}_j^w, \mathbf{w}), \mathbf{h}^*)] =$$

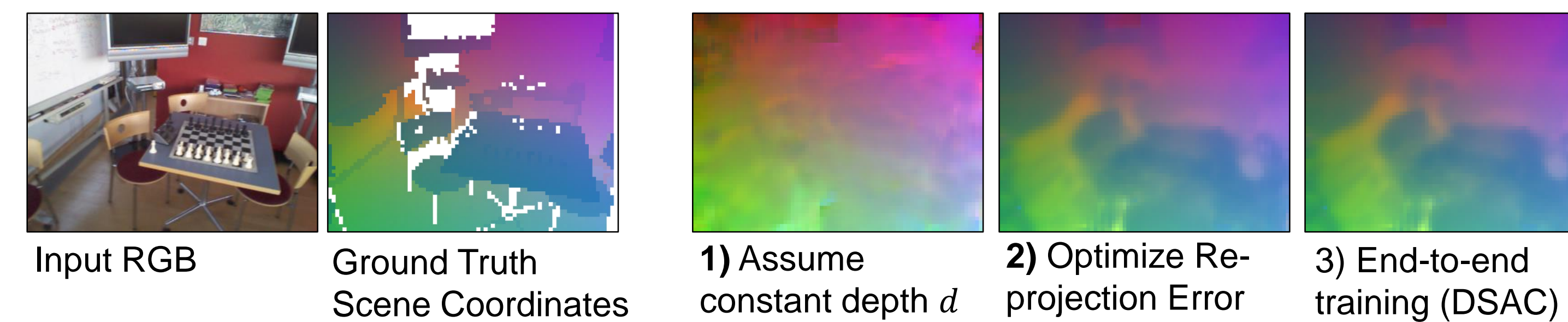
$$\mathbb{E}_{j \sim P(j|\mathbf{w})} [\ell(\cdot) \frac{\partial}{\partial \mathbf{w}} \log P(j|\mathbf{w}) + \frac{\partial}{\partial \mathbf{w}} \ell(\cdot)]$$

Updated Pipeline



Learning without a 3D Scene Model

Training scene coordinate regression in 3 stages:



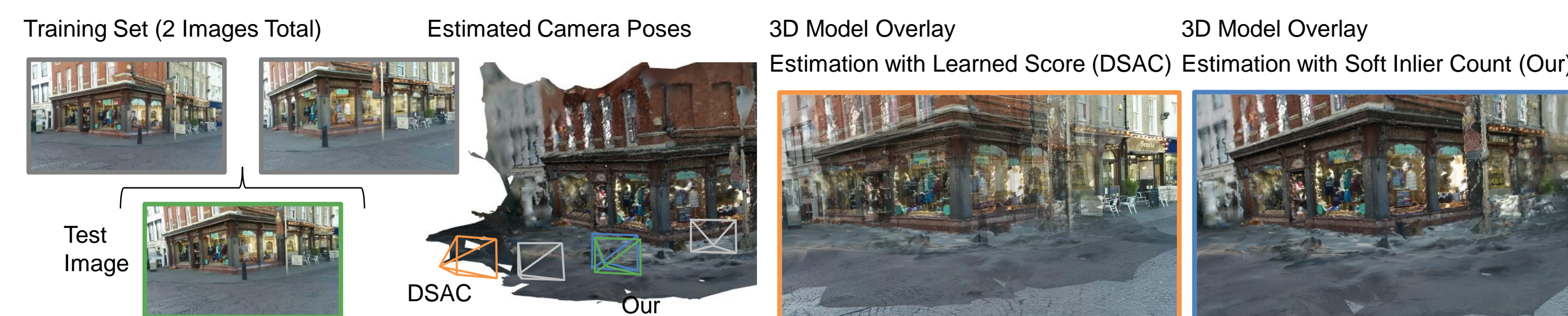
$$1) \min \sum_i \| \mathbf{y}_{i(\mathbf{w})} - \mathbf{y}_i^* \|, \text{ with } \mathbf{y}_i^* = \mathbf{h}^* \left[\frac{dx_i}{f}, \frac{dy_i}{f}, d, 1 \right]^T$$

$$2) \min \sum_i \| C\mathbf{h}^{*-1} \mathbf{y}_{i(\mathbf{w})} - \mathbf{p}_i \|$$

Hypothesis Score: Soft Inlier Count

Previously: learned $s(\mathbf{h})$ - **hard to regularize, overfits**
 Inlier Count: $s(\mathbf{h}) = \sum_i \mathbb{1}[\tau - r_i(\mathbf{h}, \mathbf{w})]$ - **not differentiable**
 Soft In. Count: $s(\mathbf{h}) = \sum_i \text{sig}(\tau - \beta r_i(\mathbf{h}, \mathbf{w}))$ - **differentiable**

$$r_i(\mathbf{h}, \mathbf{w}) = \| C\mathbf{h}^{-1} \mathbf{y}_{i(\mathbf{w})} - \mathbf{p}_i \|$$



[Sho13] "Scene Coordinate Regression Forests for Camera Relocalization in RGB-D Images", Shotton et al., CVPR'13
 [Ken15] "PoseNet: A Convolutional Network for Real-Time 6-DOF Camera Relocalization" Kendall et al., ICCV'15
 [Ken17] "Geometric Loss Functions for Camera Pose Regression with Deep Learning" Kendall and Cipolla, CVPR 2017
 [Sat16] "Efficient & Effective Prioritized Matching for Large-Scale Image-Based Localization" Sattler et al., PAMI 2016
 [Bra17] "DSAC - Differentiable RANSAC for Camera Localization", Brachmann et al., CVPR'17

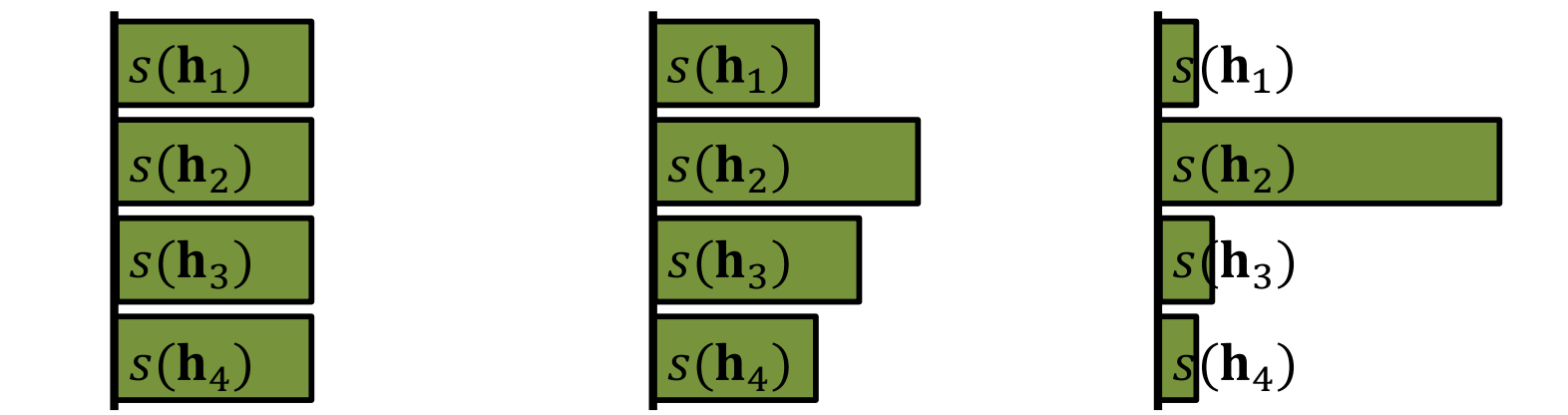
Hypothesis Score: Entropy Control

Hypothesis Distribution:

$$P(j|\mathbf{w}, \alpha) = \frac{\exp(\alpha s(\mathbf{h}_j, \mathbf{w}))}{\sum_k \exp(\alpha s(\mathbf{h}_k, \mathbf{w}))}$$

Entropy:

$$S(\alpha) = - \sum_j P(j|\mathbf{w}, \alpha) \log P(j|\mathbf{w}, \alpha)$$



Keep target entropy S^* during training by adjusting α : $\text{argmin}_\alpha |S(\alpha) - S^*|$

Differentiable Refinement

Refinement \mathbf{R} optimizes re-projection errors \mathbf{r}_j of inlier set \mathcal{J} :

$$\mathbf{R}(\mathbf{h}) = \text{argmin}_{\mathbf{h}'} \| \mathbf{r}_j(\mathbf{h}', \mathbf{w}) \|^2$$

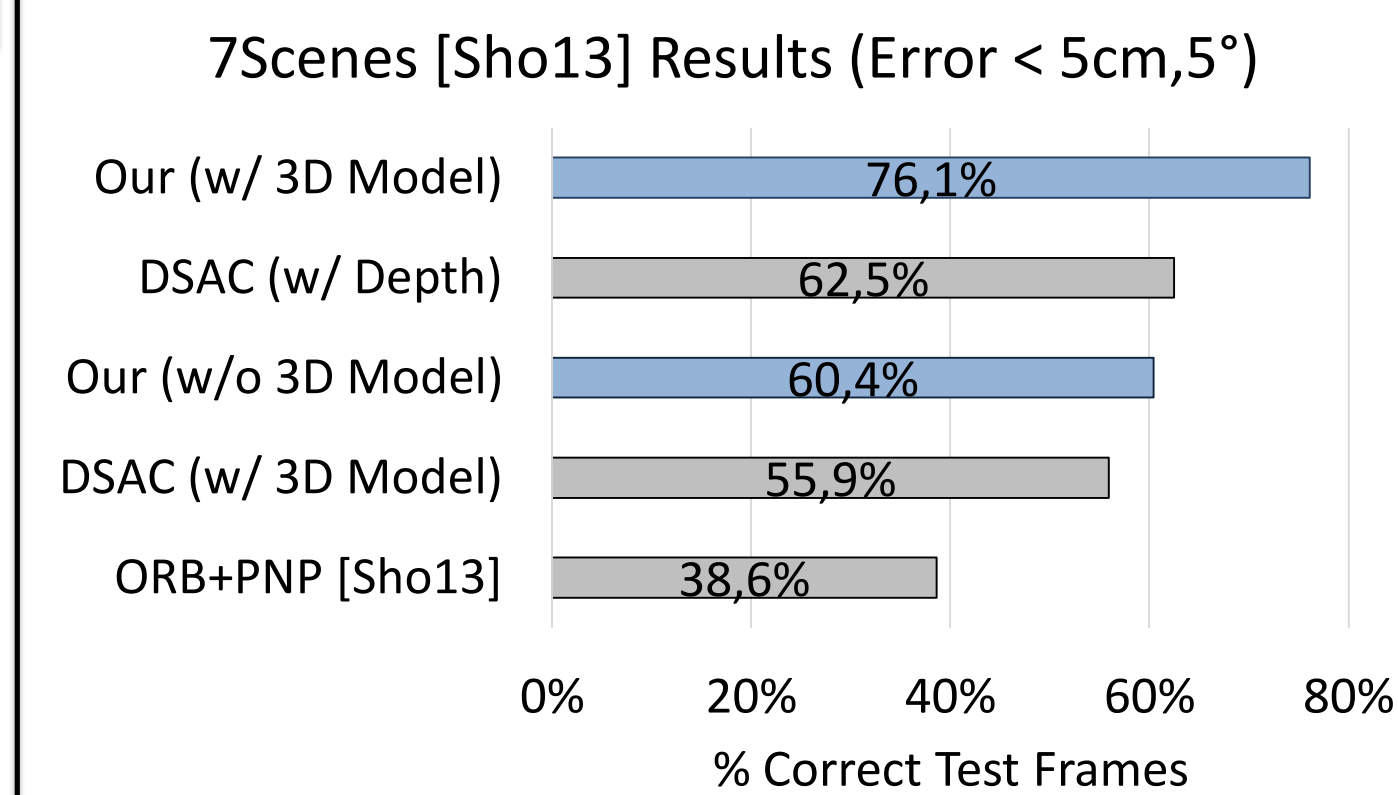
Gauss-Newton update step:

$$\mathbf{R}^{t+1} = \mathbf{R}^t - (J_{\mathbf{r}}^T J_{\mathbf{r}})^{-1} J_{\mathbf{r}}^T \mathbf{r}_j(\mathbf{R}^t, \mathbf{w})$$

Last update: $\mathbf{R}(\mathbf{h}) = \mathbf{h}_0 - (J_{\mathbf{r}}^T J_{\mathbf{r}})^{-1} J_{\mathbf{r}}^T \mathbf{r}_j(\mathbf{h}_0, \mathbf{w})$, with $\mathbf{h}_0 = \mathbf{R}^{t=\infty}(\mathbf{h})$

$$\text{Gradient approximation: } \frac{\partial}{\partial \mathbf{w}} \mathbf{R}(\mathbf{h}) \approx - (J_{\mathbf{r}}^T J_{\mathbf{r}})^{-1} J_{\mathbf{r}}^T \frac{\partial}{\partial \mathbf{w}} \mathbf{r}_j(\mathbf{h}_0, \mathbf{w})$$

Results



Cambridge Landmarks [Ken15] Results

Avg. Median Err.	w/ 3D Model	w/o 3D Model
PoseNet [Ken17]	1.43m, 2.9°	1.63m, 2.8°
Active Search [Sat16]	0.29m, 0.6°	-
DSAC [Bra17]	0.31m, 0.8°	-
Our	0.14m, 0.3°	0.19m, 0.5°

