

Coarse-to-fine Planar Regularization for Dense Monocular Depth Estimation

Stephan Liwicki Christopher Zach Ondrej Miksik Philip H. S. Torr

1 Introduction

Simultaneous localization and mapping (SLAM) using the whole image data is an appealing framework to address shortcoming of sparse feature-based methods – in particular frequent failures in textureless environments. Hence, direct methods bypassing the need of feature extraction and matching became recently popular. Many of these methods operate by alternating between pose estimation and computing (semi-)dense depth maps, and are therefore not fully exploiting the advantages of joint optimization with respect to depth and pose.

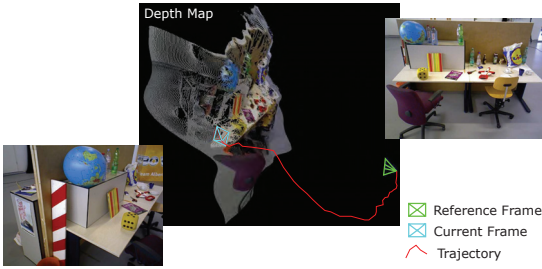


Figure 1: During keyframe-to-frame comparison a dense depth map is built. An incremental implementation allows for longer tracks with a single reference frame in real-time.

> Contributions

- >>> Global energy for planar inverse depth that is optimized iteratively
- >>> Coarse-to-fine strategy that refines depth and pose truly simultaneously
- >>> Semi-dense version, computing depth twice as fast as LSD-SLAM [1] on CPU
- >>> Quantitative pose and depth evaluation on the TUM dataset [4]

2 Global Energy

The global energy (using a single reference image) is formulated as

$$E_{All}(\mathcal{S}, \Xi) = \sum_{t=1}^T E_{Match}^{(t)}(\mathcal{S}, \xi_t) + E_{Smooth}(\mathcal{S}), \quad (1)$$

- >>> Photometric error $E_{Match}^{(t)}$
- >>> Spatial smoothing E_{Smooth} (over inverse depth values)
- >>> Camera poses $\Xi = (\xi_t)_{t=1}^T$
- >>> Planes $\mathcal{S} = (s_i^T)_{i=1}^{|\mathcal{X}|}$ (over-parametrization of $\mathcal{D} = (d_i)_{i=1}^{|\mathcal{X}|}$, where $d_i = s_i^T x_i$)

> Temporal Smoothing

We modify E_{All} to be suitable for an incremental approach for current pose ξ_T :

$$E_{All}(\mathcal{S}, \xi_T) = E_{History}^{(T)}(\mathcal{S}) + E_{Match}^{(T)}(\mathcal{S}, \xi_T) + E_{Smooth}(\mathcal{S}). \quad (2)$$

- >>> Second order approximation of $E_{History}^{(T)}$ provides temporal smoothness

> Photometric Matching Cost

We express the photometric error by

$$E_{Match}^{(T)}(\mathcal{S}, \xi_T) = \sum_{x_i \in \mathcal{X}} \|I(x_i) - I_T(W(x_i, d_i, \xi_T))\|_{\tau_{Match}}. \quad (3)$$

- >>> Warping $x_i' = W(x_i, d_i, \xi_t) = \text{hom}(\mathbf{R}_i^T(x_i - \mathbf{t}_i d_i)) = \text{hom}(\mathbf{R}_i^T(x - \mathbf{t}_i s_i^T x))$
- >>> Robust smooth truncated quadratic $\|\cdot\|_{\tau_{Match}}$ given in [2]

> Local Spatial Plane Regularizer

The smoothness constraint $E_{Smooth}(\mathcal{S})$ is expressed in inverse depth:

$$E_{Smooth}(\mathcal{S}) = \lambda_{Smooth} \sum_{x_i \in \mathcal{X}} \sum_{x_j \in \mathcal{N}_i} \|s_i^T x_i - s_j^T x_j\|_{\tau_{Smooth}}. \quad (4)$$

- >>> Balancing term λ_{Smooth}
- >>> Smoothness assumption inspired by stereo setups

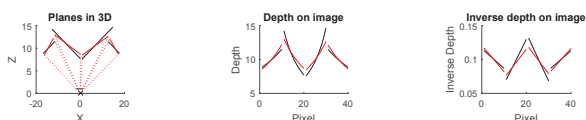


Figure 2: Top view of planes. Pixels in 3D space are aligned via smoothing in the inverse depth image.

3 Optimization Strategy

We find pose as well as depth through Levenberg-Marquardt optimization.

- >>> Locally linearize I_T (as in KLT)
- >>> Graduated optimization [3]
 - > scale-space pyramid
 - > proposed restricted depth



Figure 3: We incrementally increase depth resolution.

Input: Keyframe I and images $(I_i)_{i=1}^L$.
Output: Final pose ξ and depth hypothesis \mathcal{S} .

- 1: $s_i \leftarrow [0 \ 0 \ 1]^T$ and initial depth certainty $\lambda_i \leftarrow 0$ for all $x_i \in \mathcal{X}$.
- 2: compute resolution pyramid for the keyframe I .
- 3: $\xi \leftarrow (\mathbf{I} \in \mathbb{R}^{3 \times 4}, [0 \ 0 \ 0]^T)$
- 4: for each frame I_i do
- 5: compute resolution pyramid for the frame I_i .
- 6: for each pyramid level do
- 7: optimize ξ via lie algebra $so(3)$ through Levenberg-Marquardt.
- 8: repeat
- 9: update ξ (and $s_i \leftarrow s_i + \mathbb{I}_c(x_i)\Delta$, if applicable).
- 10: introduce new component Δ_c .
- 11: estimate $\mathbb{I}_c(x_i)$ via eigenvector of $\sum_{x_i \in \mathcal{X}} \nabla_x \nabla_x^T$.
- 12: optimize ξ and Δ_c through Levenberg-Marquardt.
- 13: until improvement below $\epsilon_{Complete}$ or maximum C reached
- 14: end for
- 15: update precision λ_i and depth s_i^* for temporal constraint.
- 16: end for

Algorithm 1: Dense Incremental Planar Depth Estimation

> Coarse-to-fine Depth Perception

We enforce low complexity in the update of \mathcal{S} :

$$s_i = s_i^* + \sum_{c=1}^C \mathbb{I}_c(x_i) \Delta_c, \quad (5)$$

- >>> Indicator function $\mathbb{I}_c : \mathcal{X} \rightarrow \{+1, -1\}$ (updates are constrained to 2^C values)
- >>> Greedily add one component Δ_c at a time (3 (Δ_c) + 6 (SE(3)) parameters)

> Advantages in Optimization Conditions

- >>> Regularization across all pixels is enforced for the update of planes
- >>> The image planes encode inverse depth hierarchically from coarse to fine
- >>> The incremental algorithm enables fast, joint optimization of pose and depth
- >>> Depth updates are simple, but yield rich depth maps after repeated updates

4 Evaluation

We evaluate our methods, dense DIP and semi-dense SIP, on 6 TUM and 7 own videos.

> Quantitative Results

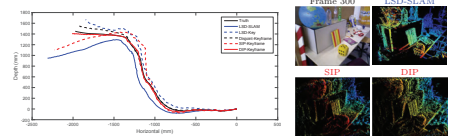


Figure 4: Quantitative trajectory (left) and its depth (right) of an example sequence.

- >>> DIP (one keyframe) performs similar to LSD-SLAM (many keyframes) for pose
- >>> DIP benefits from larger baselines for depth estimation
- >>> LSD-Key (LSD-SLAM with single reference frame) is less favourable
- >>> Our semi-dense SIP performs well for small motion
- >>> The disjoint (alternating) version is worse in virtually all experiments

> Qualitative Results

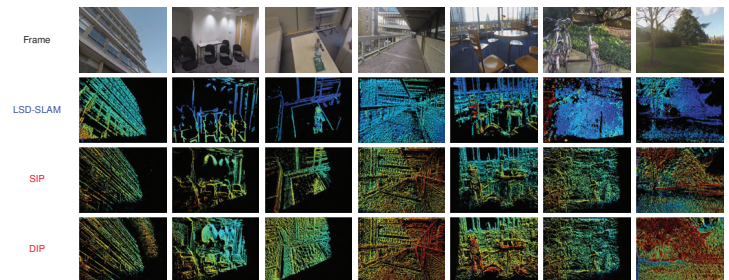


Figure 5: Qualitative depth results for our videos.

- >>> SIP, DIP produce more globally consistent depths in contrast to LSD-SLAM
- >>> Even in non-planar scenes, our method produces sensible results
- >>> Our method converges to local minima if initial $s_i \leftarrow [0 \ 0 \ 1]^T$ is significantly wrong

> Running time

- >>> DIP is highly parallel and computes in real-time on GPU
- >>> SIP estimates depth and pose twice as fast as LSD-SLAM using CPU (30 fps)

[1] J. Engel, T. Schöps, and D. Cremers. LSD-SLAM: Large-Scale Direct Monocular SLAM. In *ECCV'14*, pages 834 – 849, 2014.
 [2] H. Li, R. Sumner, and M. Pauly. Global Correspondence Optimization for Non-Rigid Registration of Depth Scans. *European Symp. Geometry Processing*, 27(6):1421–1430, 2009.
 [3] H. Mobahi and J. Fisher. On the Link between Gaussian Homotopy Continuation and Convex Envelopes. In *Int. Conf. Energy Minimization Methods Computer Vision and Pattern Recognition, EMCCV'15*, pages 43–56, 2015.
 [4] J. Sturm, N. Engelhard, F. Endres, W. Burgard, and D. Cremers. A Benchmark for the Evaluation of RGB-D SLAM Systems. In *ROS'12*, 2012.