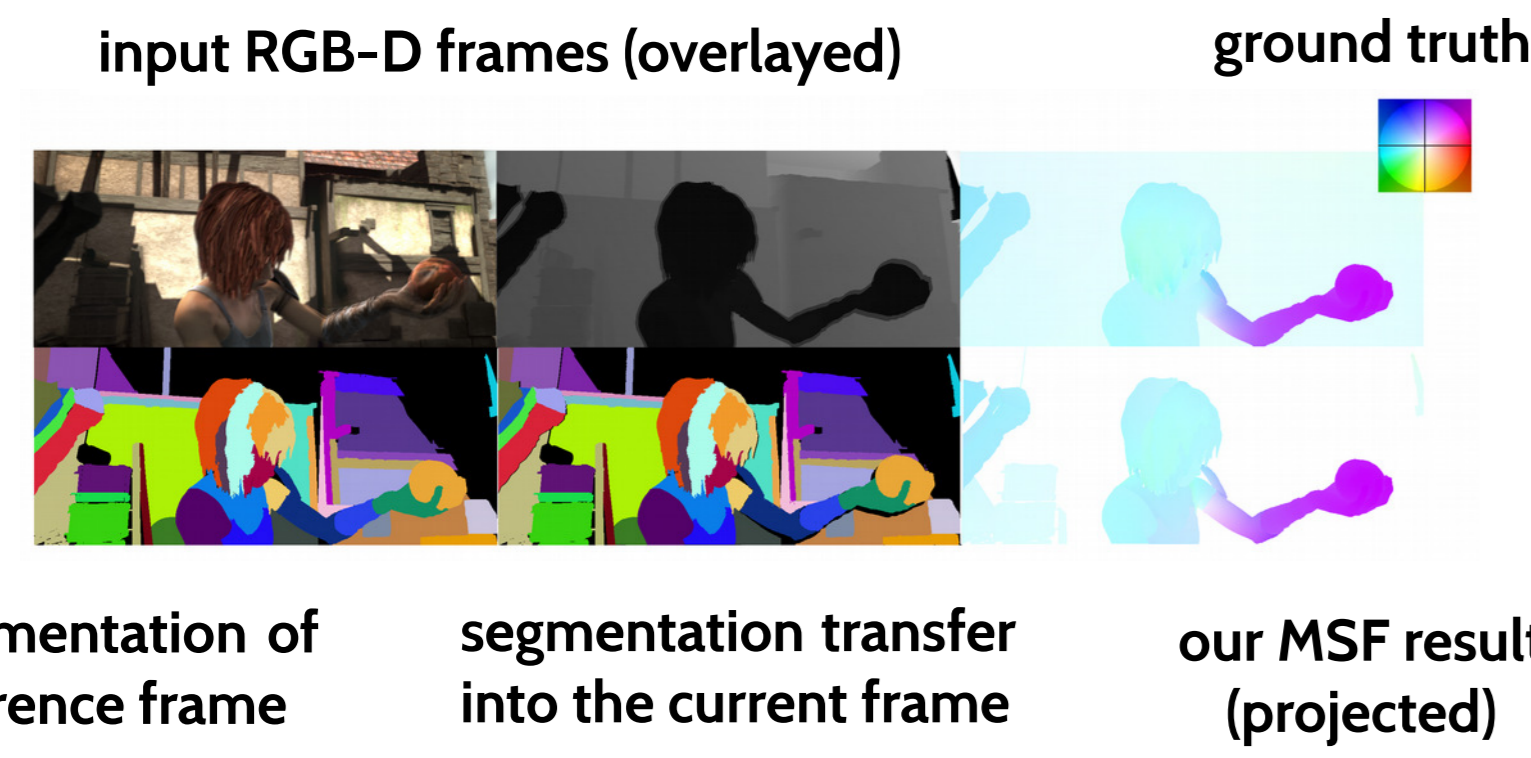


Multiframe Scene Flow with Piecewise Rigid Motion

Vladislav Golyanik, Kihwan Kim, Robert Maier, Matthias Nießner, Didier Stricker and Jan Kautz

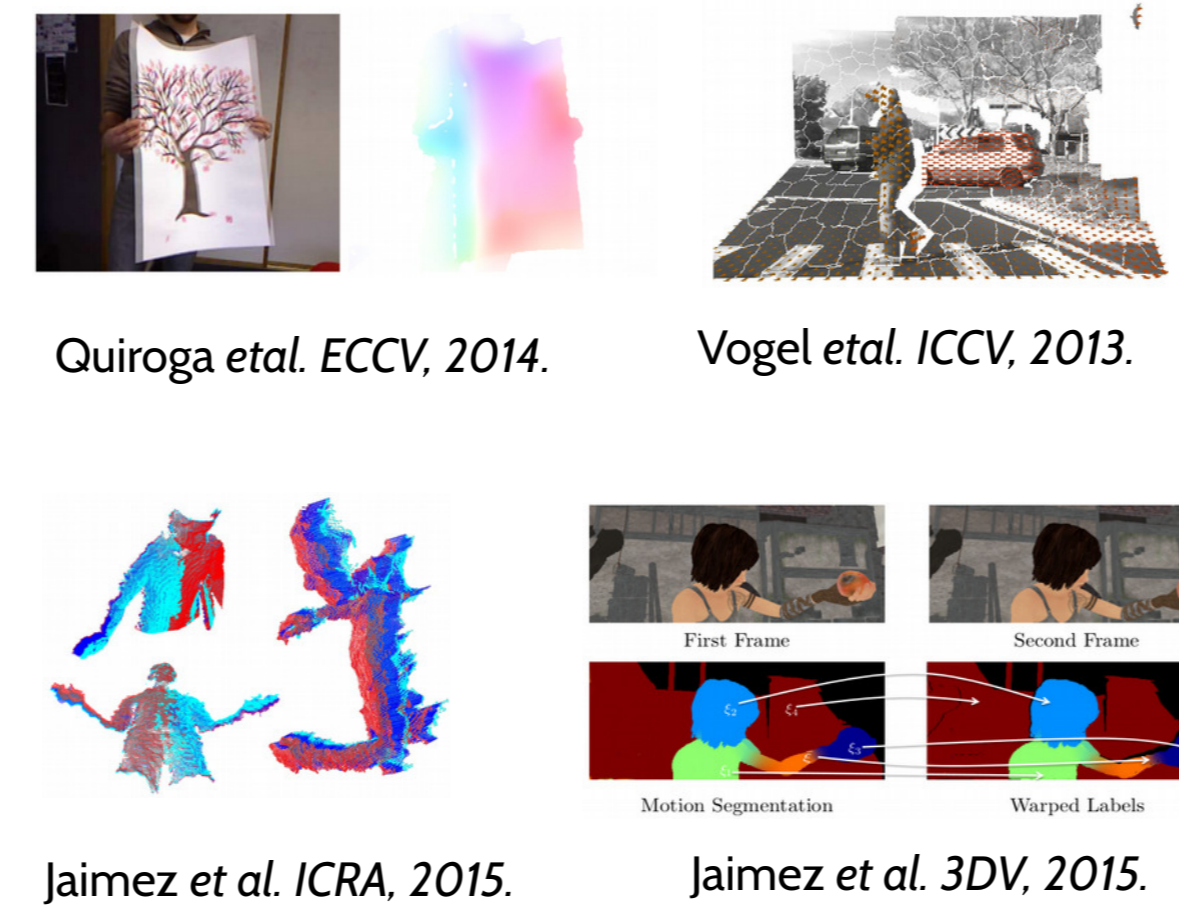
Overview



Contributions

- transformations are parametrized by **piecewise rigid motions**
- **depth channel** is used to obtain oversegmentation of the scene
- segmentation of a scene is kept **fixed**
- a global scene-flow formulation over **multiple frames**
- take advantage of **point set registration (projective point-to-plane ICP term)**
- **lifting function** for coherent segment transformations
- using the efficient framework of non-linear least squares

Related work



Proposed Energy Functional

$$\mathcal{E}(\mathbf{T}^1, \mathbf{T}^2, \dots, \mathbf{T}^{|Z|}, \mathbf{w}) = \sum_{\zeta \in Z} \alpha_{\zeta} \mathcal{E}_{\text{data}}(\mathbf{T}^{\zeta}) + \sum_{\zeta \in Z} \beta_{\zeta} \mathcal{E}_{\text{pICP}}(\mathbf{T}^{\zeta}) + \gamma_{\zeta} \sum_{\zeta \in Z} \mathcal{E}_{\text{l.reg.}}(\mathbf{T}^{\zeta}, \mathbf{w}) + \eta \mathcal{E}_{\text{r.opt.}}(\mathbf{w}) + \sum_{\zeta=3}^{|Z|} \lambda_{\zeta} \mathcal{E}_{\text{c.}}(\mathbf{T}^{\zeta})$$

data term (brightness constancy) + projective ICP term (point-to-plane) + lifted segment pose regularizer + robust weight optimizer + multiframe pose concatenation term

I $\mathcal{E}_{\text{data}}(\mathbf{T}) = \sum_k \sum_{\mathbf{p} \in \Omega_k} \|(\mathcal{I}_1(\mathbf{p}) - \mathcal{I}_2(\pi(g(\mathbf{T}_k, \pi^{-1}(\mathbf{p}, z_p))))\|_{\epsilon}$

II $\mathcal{E}_{\text{pICP}}(\mathbf{T}) = \sum_k \sum_{\mathbf{p} \in \Omega_k} \|((g(\mathbf{T}_k, \mathbf{p}) - \mathbf{p}^{\text{corr}}) \cdot \mathbf{n}_p)\|_{\epsilon}$

III $\mathcal{E}_{\text{c.}}(\mathbf{T}^{1,m}) = \sum_k \|(\mathbf{T}_k - \mathbf{T}_k^{m-1,m} \dots \mathbf{T}_k^{l,l+1})\|_2$

Huber norm: $\|a\|_{\epsilon} = \begin{cases} \frac{1}{2}a^2, & \text{for } |a| \leq \epsilon \\ \epsilon(|a| - \frac{1}{2}\epsilon), & \text{otherwise} \end{cases}$

global optimization over multiple frames

References

[1] R. Achanta et al. Slic superpixels compared to state-of-the-art superpixel methods. T-PAMI, 2012.
 [2] P. F. Felzenszwalb and D. P. Huttenlocher. Efficient graph-based image segmentation. IJCV, 2004.
 [3] A. Gaidon, Q. Wang, Y. Cabon, and E. Vig. Virtual worlds as proxy for multi-object tracking analysis. In CVPR, 2016.
 [4] M. Jaimez, M. Souli, J. González-Jiménez, and D. Cremers. A primal-dual framework for real-time dense rgb-d scene flow. In ICRA, 2015.
 [5] D. J. Butler, J. Wulff, G. B. Stanley, and M. J. Black. A naturalistic open source movie for optical flow evaluation. In ECCV, 2012.
 [6] J. Quiroga, T. Brox, F. Deynmay, and J. L. Crowley. Dense semi-rigid scene flow estimation from RGB-D images. In ECCV, 2014.
 [7] J. Stueckler and S. Behnke. Efficient dense rigid-body motion segmentation and estimation in rgb-d video. IJCV, 2015.
 [8] B. Taetz, G. Bleser, V. Golyanik, and D. Stricker. Occlusion-aware video registration for highly non-rigid objects. In WACV, 2016.
 [9] C. Vogel, K. Schindler, and S. Roth. Piecewise rigid scene flow. In ICCV, 2013.
 [10] C. Zach, T. Pock, and H. Bischof. A duality based approach for realtime tv-l optical flow. In GCPR, 2007.

Energy Optimization

We minimize the energy using **Levenberg-Marquardt (ceres solver, C++, multithreaded)**:

$$\mathbf{x}' = \arg \min_{\mathbf{x}} \|\mathbf{F}(\mathbf{x})\|_2^2 \quad \mathbf{F}(\mathbf{x} + \Delta\mathbf{x}) \approx \mathbf{F}(\mathbf{x}) + \mathbf{J}(\mathbf{x})\Delta\mathbf{x} \quad (\mathbf{J}(\mathbf{x})^T \mathbf{J}(\mathbf{x}) + \lambda \text{diag}(\mathbf{J}(\mathbf{x})^T \mathbf{J}(\mathbf{x}))) \Delta\mathbf{x} = -\mathbf{J}(\mathbf{x})^T \mathbf{F}(\mathbf{x})$$

The total number of residuals:

$$M = {}^N C_2 (n_c + n_D) + ({}^N C_2 + 1) n_{pp} + n_c$$

of individual frame-to-frame (two frame) combinations # of valid pixel pairs (two frames) # of valid 3D point pairs (two frames) # of weights (lifting) # of concatenated transformations $n_c = K(N - 2)$

of pose pairs (number of non-zero elements in the segment adjacency matrix), for every two frame case $n_{pp} = {}^N C_2 \sum \psi_{k,l}$

Experiments

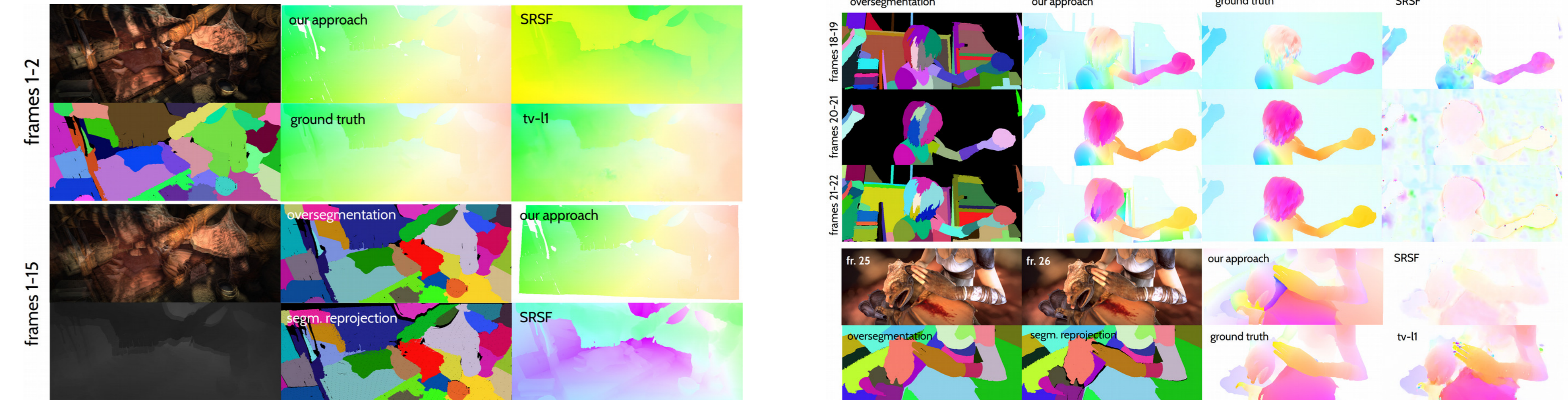
In the experimental evaluation we use:

- MPI SINTEL [5]
- virtual KITTI [3]
- Bonn multibody data set [7]
- own RGB-D recordings

... and compare the following methods:

- Primal-Dual Flow [4]
- Semi-Rigid Scene Flow [6]
- Multi-Frame Optical Flow [8]
- tv-l1 optical flow [10]

End Point Error is defined as $\|(u - u_{GT}), (v - v_{GT})\|$
 $(u, v)^T$ is a projected flow vector
 $(u_{GT}, v_{GT})^T$ is a ground truth vector

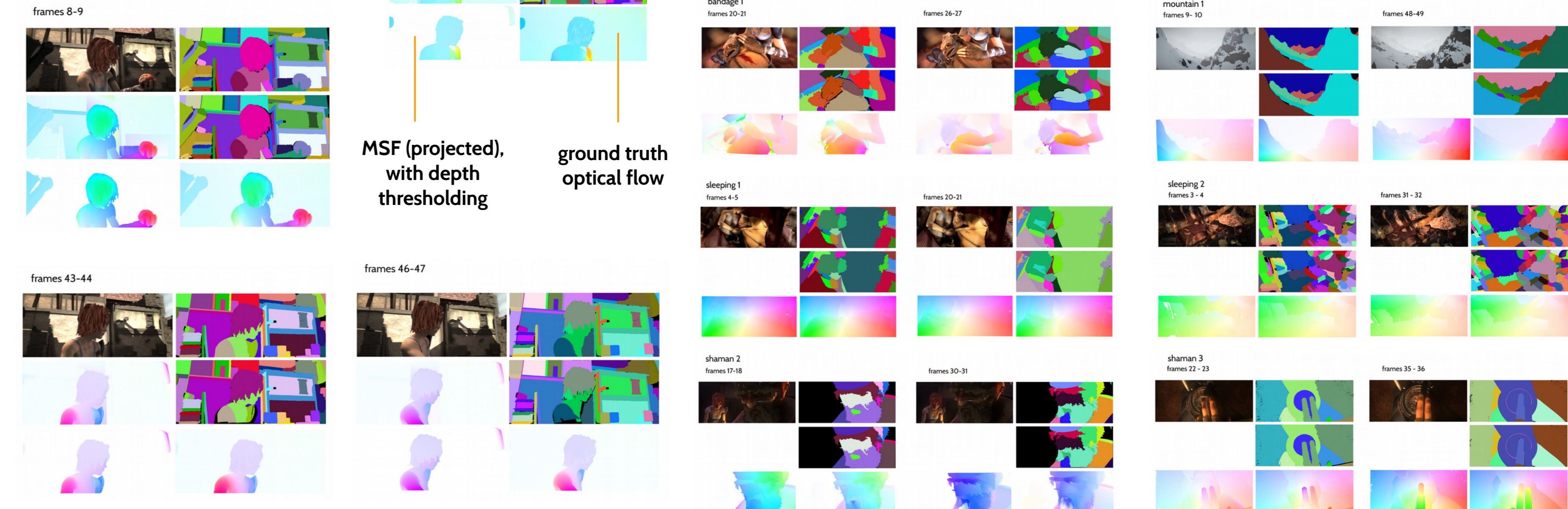


	alley1	bandage1	sleeping2_rigid
SRSF [6]	2.46122/2.40833	2.47801/2.46389	1.13384
MSF	0.740127	1.69865	0.307526

comparisons of average EPE between scene flow projections and the ground truth optical flow on the MPI SINTEL [5]

	2 fr.	3 fr.	4 fr.	2 fr. (r)	1. d. (r)
SRSF [6]	274	n.a.	n.a.	87.5	84.5
MSF	49	221	541	90	254

runtime comparisons of SRSR [6] and the proposed MSF for several configurations



Results of several RGB-D scene flow and optical flow approaches on the Bonn multibody data set [7] and an example of an urbane driving scene processing from the vKITTI data set [3] by MSF (bottom right)

Segmentation transfer from the reference frame to three other frames in *alley1*

Segmentation transfer on the Bonn watering can sequence with Felzenszwalb segmentation [2] (top row) and SLIC segmentation [1] (bottom row)

