

# Dense 3D scene flow estimation for locally rigid scenes

Konrad Schindler

Photogrammetry and Remote Sensing, ETH Zürich



# 3D Scene Flow

- “Joint stereo and optical flow”
- Given  $\geq 2$  video frames from  $\geq 2$  different viewpoints



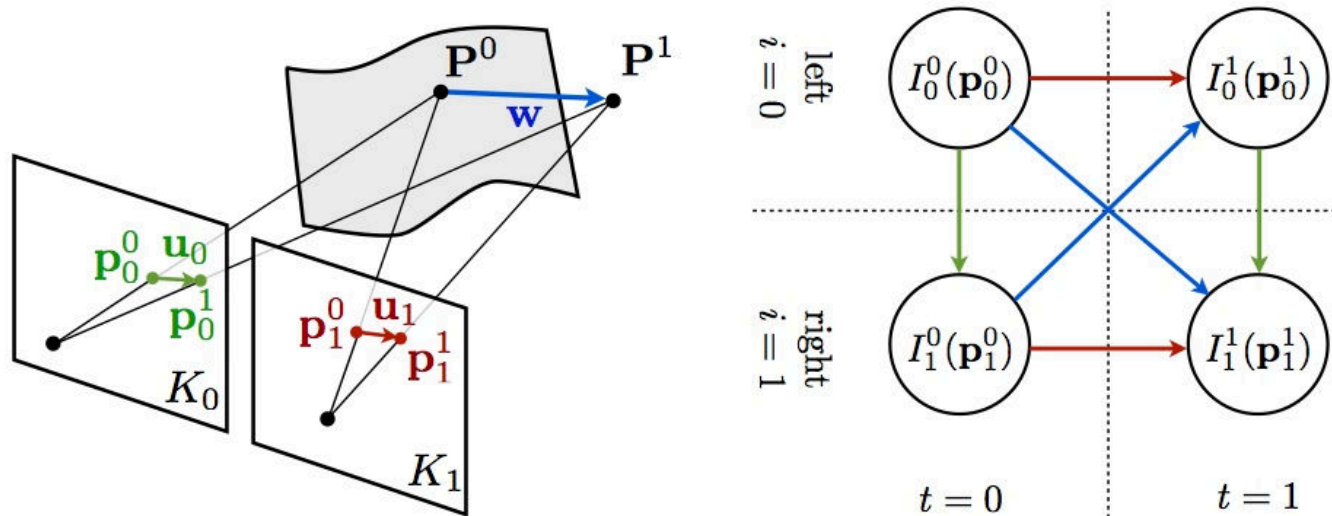
# 3D Scene Flow

- “Joint stereo and optical flow”
- Given  $\geq 2$  video frames from  $\geq 2$  different viewpoints
- Estimate dense 3D shape and 3D motion field



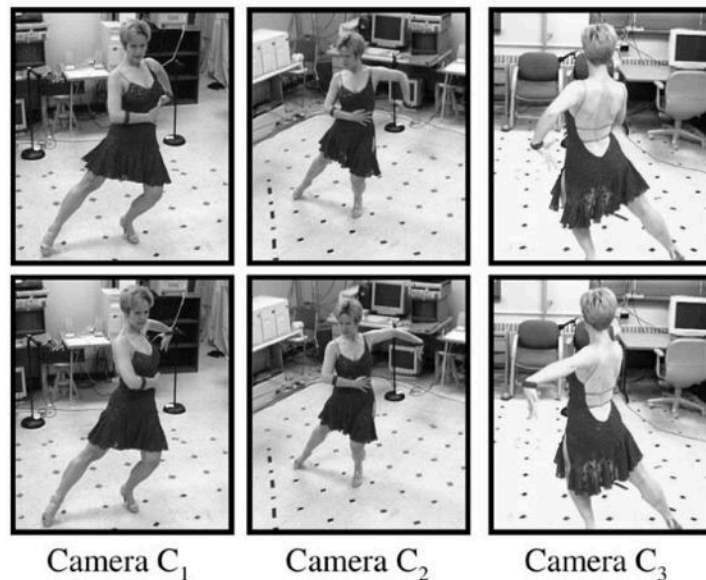
# 3D Scene Flow

- 4 unknowns per point: depth + 3D motion vector
- stereo/flow constraints between any two images
- (for simplicity, will refer to the 2-view case in the talk)

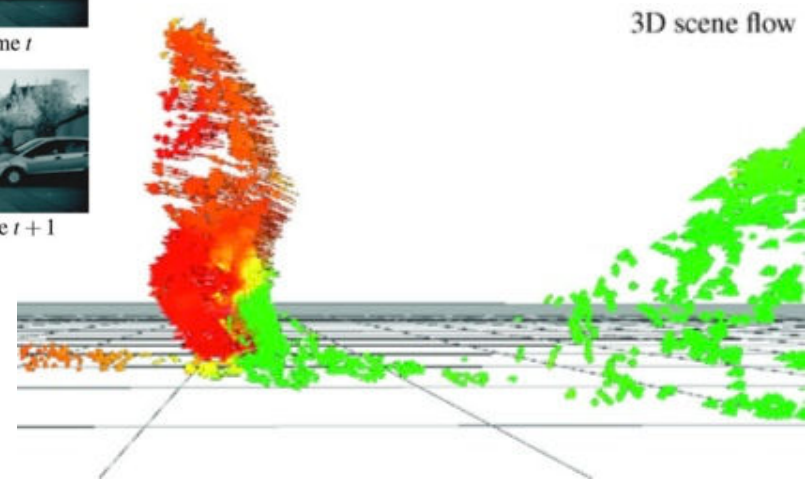


# Applications

- Entertainment, 3D TV
- Motion capture
- Driver assistance, autonomous robots



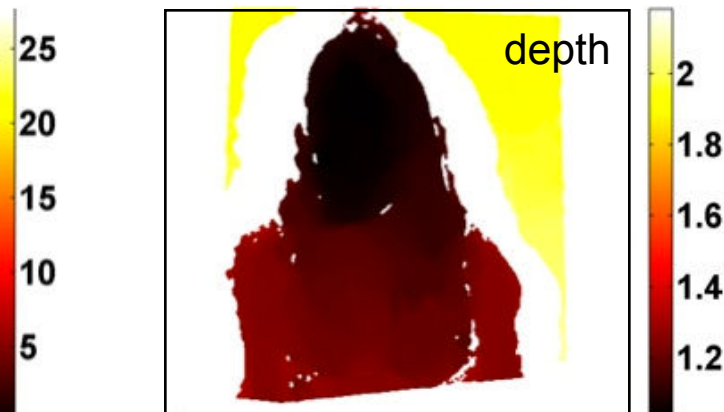
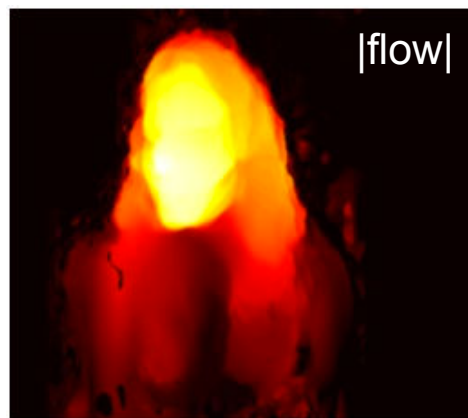
[Vedula, Baker, Rander,  
Collins, Kanade 1999]



[Wedel, Rabe, Vaudrey,  
Brox, Franke, Cremers 2008]

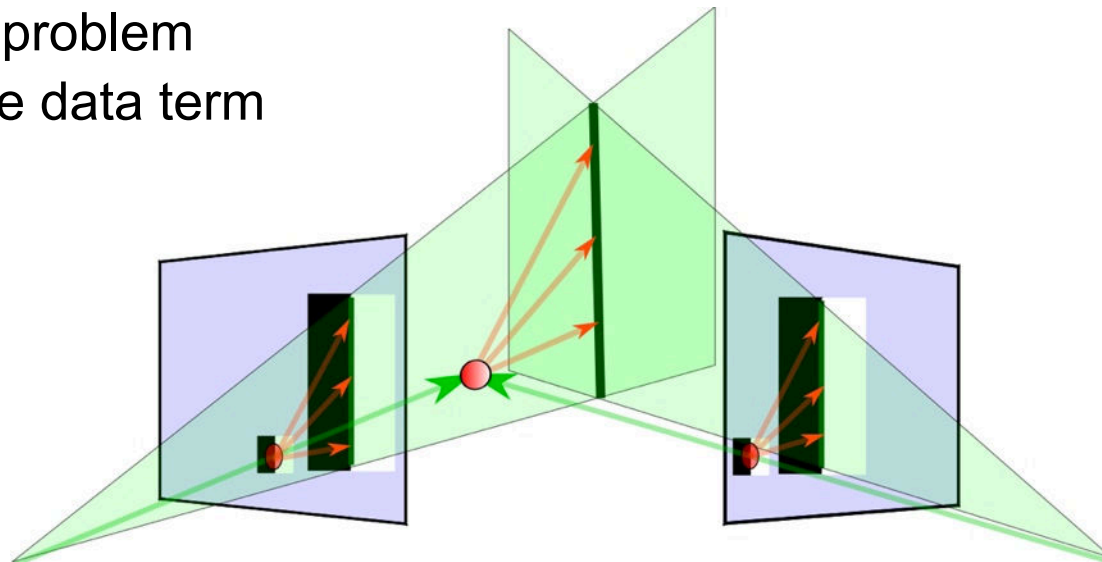
# Estimation

- First 2D flow, then stereo [Vedula et al. 1999]
- First stereo, then 2D flow [Pons, Keriven, Faugeras 2007]  
[Wedel et al. 2008]
- Everything jointly [Huguet, Devernay 2007]  
[Basha, Moses, Kiryati 2010]
- Usually parametrized w.r.t. a **reference image**



# Scene Flow is ill-posed

- Available constraint (data term): the flow field should be compatible with all observed images
  - e.g. brightness constancy, cross-correlation, Census, ...
- Like for stereo and optical flow, this is insufficient to determine the flow field
  - 3D aperture problem
  - uninformative data term



# Standard regularization

- surface / motion field is piecewise smooth

$$\hat{\mathbf{u}} = \operatorname{argmin} E(\mathbf{u}) \quad , \quad E(\mathbf{u}) = E_{data}(\mathbf{u}) + \lambda \cdot E_{smooth}(\mathbf{u})$$

- robust penalty to allow for discontinuities, e.g. total variation (TV-L1)

$$E_{data} = \int \rho(\nabla \mathbf{u}) d\mathbf{x} \quad , \quad \rho(\nabla \mathbf{u}) = |\nabla \mathbf{u}|$$

- widely used in stereo and optical flow estimation

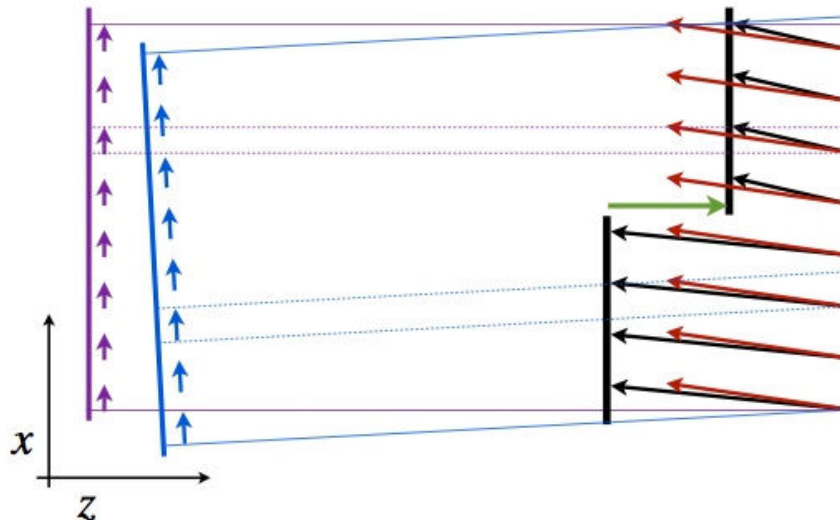


# TV-L1 applied to 3D scene flow

- Total variation of depth and motion field

$$E_S^{\text{TV}}(d, \mathbf{w}) = \int_{\Omega} \rho(\nabla d) + \rho(\nabla w_x) + \rho(\nabla w_y) + \rho(\nabla w_z) dx$$

- Does not work well for narrow baselines
- Isotropic smoothing is biased against depth discontinuities



# Rigidity

- Observation: many scenes of practical importance consist of relatively few rigid, independently moving parts
- Develop a regularizer that encourages **rigidity** rather than smoothness



# Attempt 1 - local rigidity

- For each scene point, penalize the deviation from rigidity
  - look at a small patch around the point
  - find the “non-rigid motion residual”, i.e. the difference between the observed motion and its projection onto the rigid motion subspace
  - smoothness term is a robust function of this residual

$$E_S^R(\mathbf{w}) = \int_{\Omega} \psi(v^R(\mathbf{x}; \mathbf{w})) \, d\mathbf{x}.$$

$$v^R(\mathbf{x}; \mathbf{w}) = \int_{\mathcal{C}(\mathbf{x})} \|\mathbf{r}(\mathbf{y}; \mathbf{w} |_{\mathcal{C}(\mathbf{x})}) - \mathbf{w}(\mathbf{y})\|_2^2 \eta(\mathbf{x}, \mathbf{y}) \, d\mathbf{y}$$



[Vogel, Roth, Schindler ICCV'11]

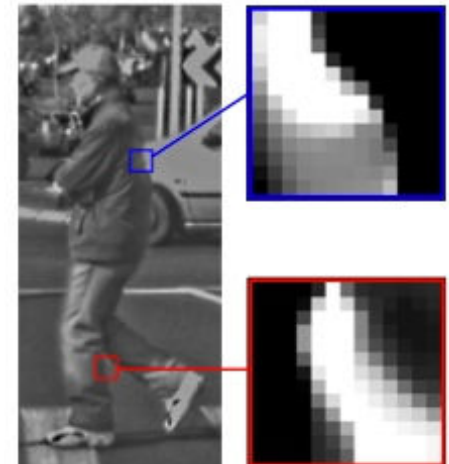
# Local rigidity

- For small motions this can be done quite efficiently
  - linearize rotation, discretize to pixel grid
  - non-rigid motion residual has linear closed form
  - it is not necessary to explicitly compute the rigid motion

$$E_S^R(\mathbf{w}) = \sum_{c \in C} \psi(v_{\text{dsc}}^R(c; \mathbf{w})) \quad R \approx I + \alpha[\mathbf{r}]_{\times}$$

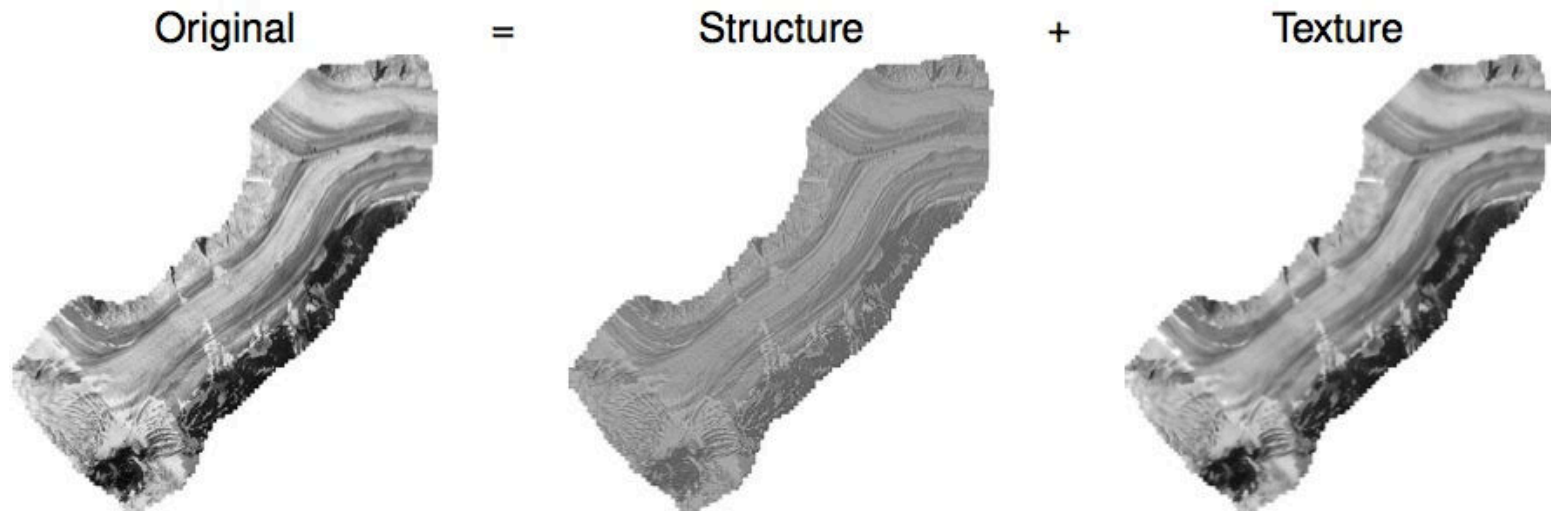
$$v_{\text{dsc}}^R(c; \mathbf{w}) = \left\| A_c \mathbf{w}_{(c)} - \mathbf{w}_{(c)} \right\|_{N_c}^2$$

- Pixels in the patch are weighted to avoid fitting across depth and motion boundaries



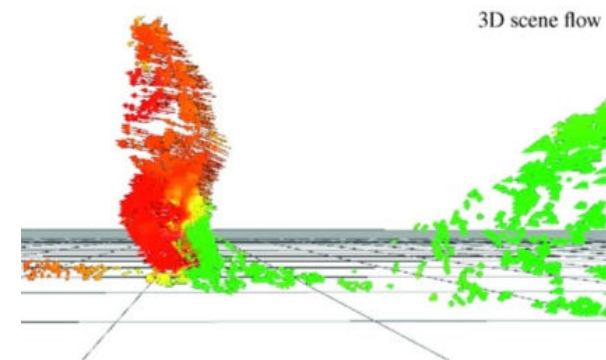
# Inference

- variational energy minimization
- usual tricks from optical flow can (and should) be used
  - auxiliary dual variables to decouple data and smoothness terms
  - course-to-fine scheme
  - careful gradient interpolation
  - structure-texture decomposition



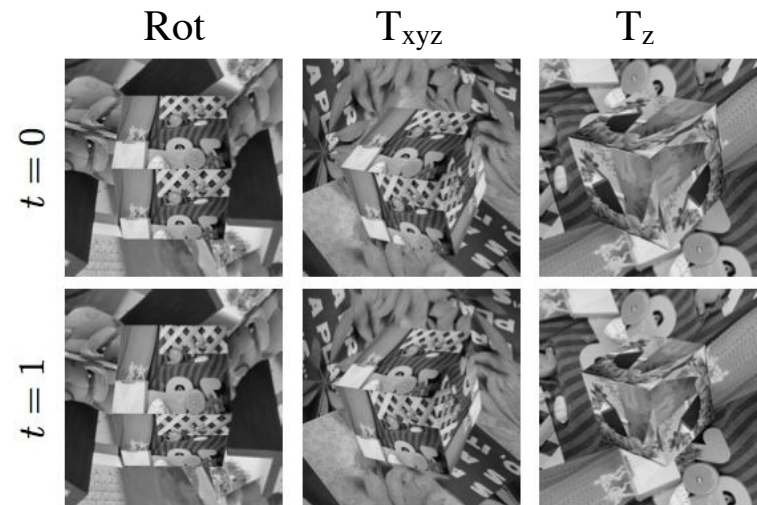
# Evaluating scene flow

- Ground truth
  - no easy way to obtain dense ground truth
  - no existing benchmark data
  - quantitative results only for synthetic scenes
  - (note, some public ground truth is wrong)
- Error measures
  - angular error of 3D flow component
  - RMS of depth, normalized to scene extent
  - RMS of flow, normalized to scene extent
  - 2D angular of reprojected 2D flow field
  - 2D endpoint error of reprojected flow field



# Results

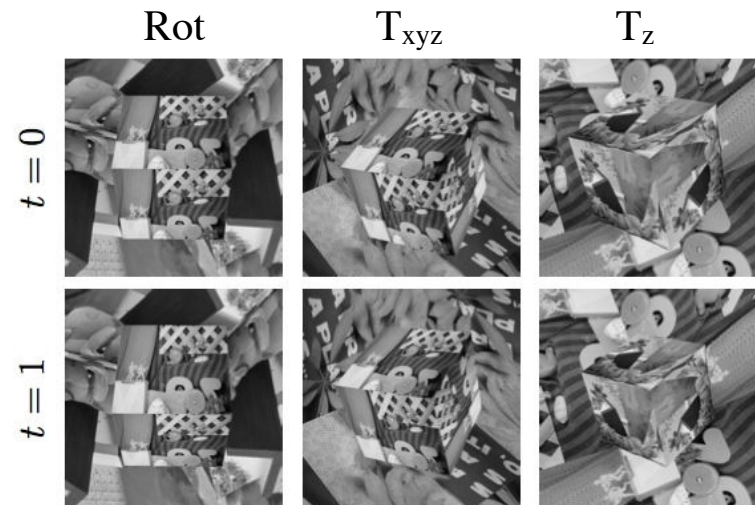
- Synthetic scenes with known ground truth



SCENE		3D ERROR			2D ERROR	
		$AAE_w$	$NRMS_w$	$NRMS_d$	$AAE$	$AEP$
Rot	Rig	<b>4.5°</b>	<b>7.3%</b>	11.7 %	<b>1.6°</b>	0.36
	TV	8.5°	9.8%	<b>11.6%</b>	1.7°	<b>0.35</b>
$T_{xyz}$	Rig	<b>2.5°</b>	<b>11.9%</b>	11.8 %	<b>1.5°</b>	<b>0.39</b>
	TV	8.6°	25.6 %	<b>11.7%</b>	2.3°	0.42
$T_z$	Rig	<b>3.9°</b>	<b>14.0%</b>	<b>9.9%</b>	<b>1.8°</b>	<b>0.35</b>
	TV	7.8°	15.3%	10.7 %	2.4°	0.37

# Results

- Synthetic scenes with known ground truth

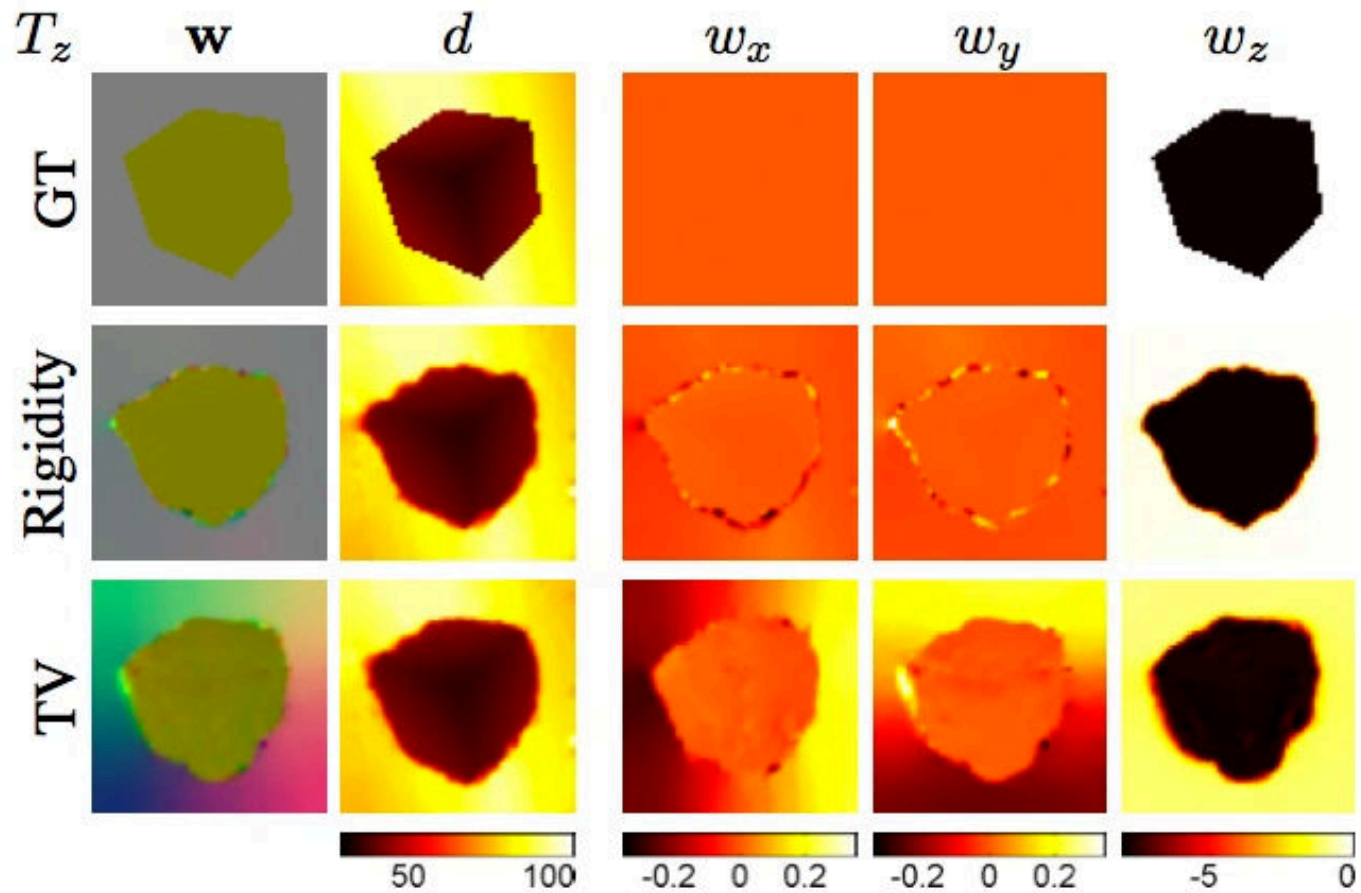


SCENE		AAE <sub>w</sub> [°]			NRMS <sub>w</sub> [%]		
		—	OC	OC&DC	—	OC	OC&DC
Rot	Rig	<b>4.5</b>	<b>4.1</b>	<b>3.3</b>	<b>7.3</b>	<b>6.7</b>	<b>4.5</b>
	TV	8.5	8.2	7.4	9.8	9.4	7.9
$T_{xyz}$	Rig	<b>2.5</b>	<b>2.1</b>	<b>1.5</b>	<b>11.9</b>	<b>10.3</b>	<b>6.4</b>
	TV	8.6	8.0	7.7	25.6	24.4	23.6
$T_z$	Rig	<b>3.9</b>	<b>2.5</b>	<b>1.2</b>	<b>14.0</b>	<b>10.6</b>	<b>5.0</b>
	TV	7.8	6.5	4.9	15.3	13.3	8.2



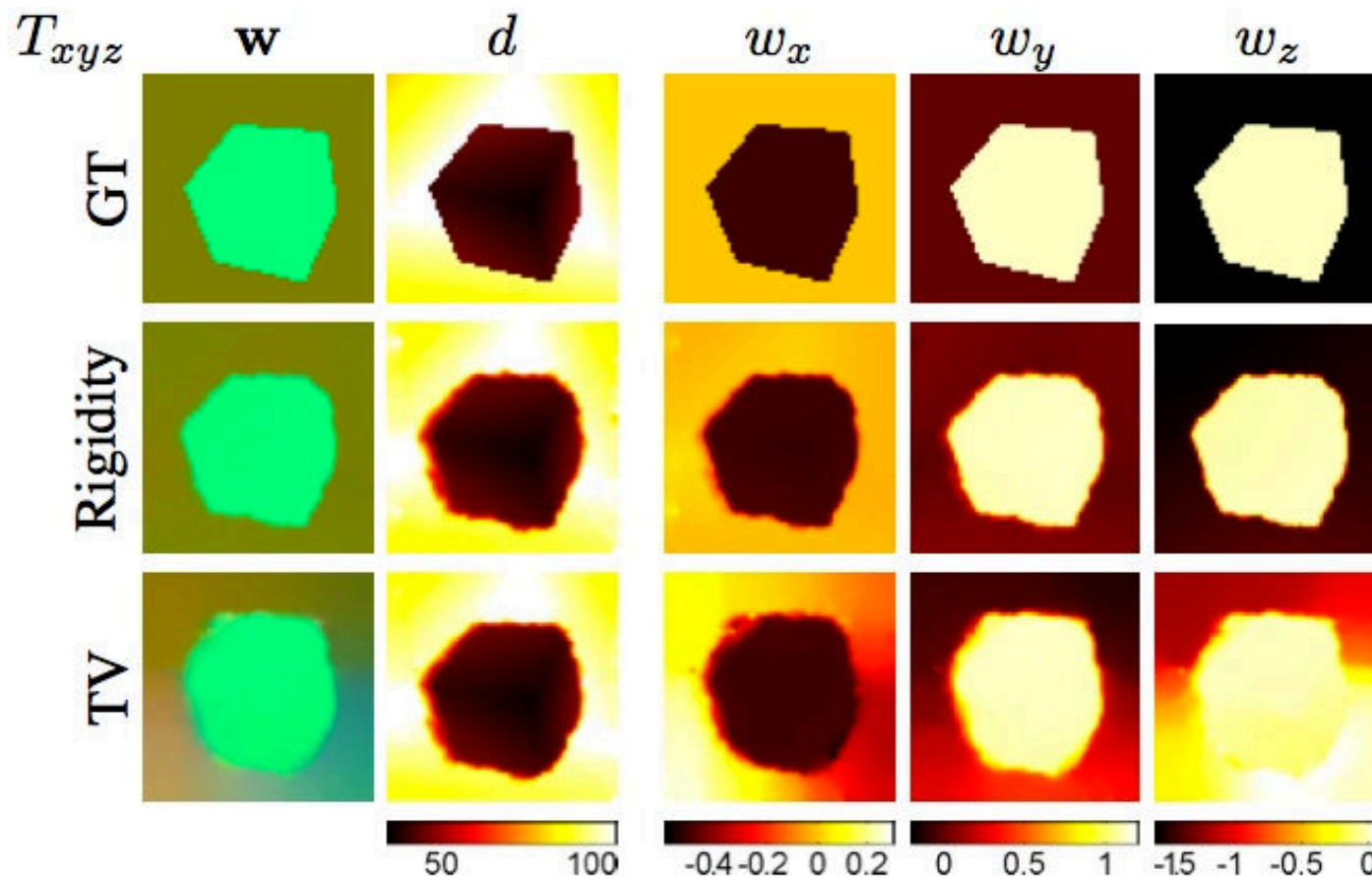
# Results

- Qualitative results



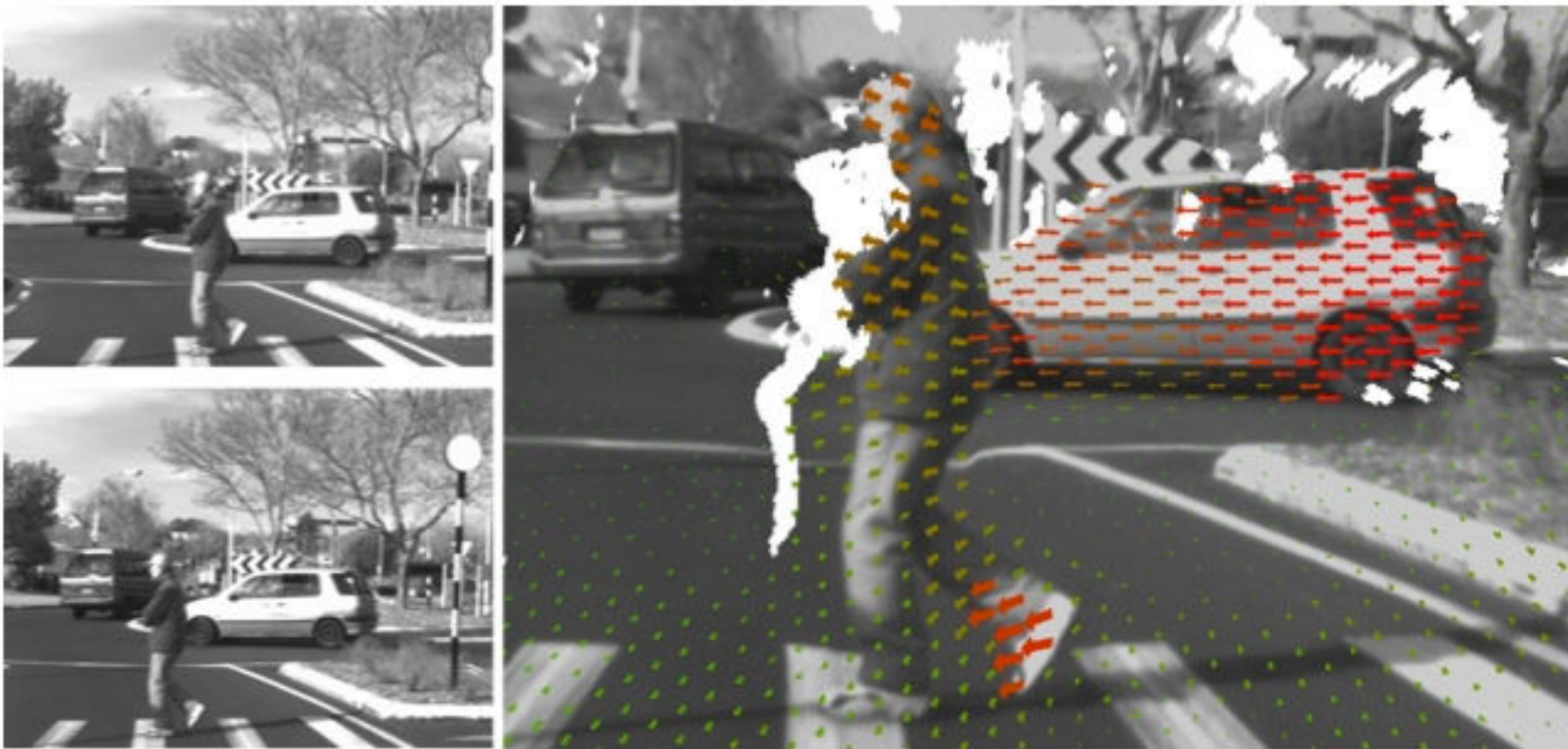
# Results

- Qualitative results



# Results

- Stereo camera on car



# Results

- Maria (three views)



## Attempt 2 - piecewise planarity and rigidity

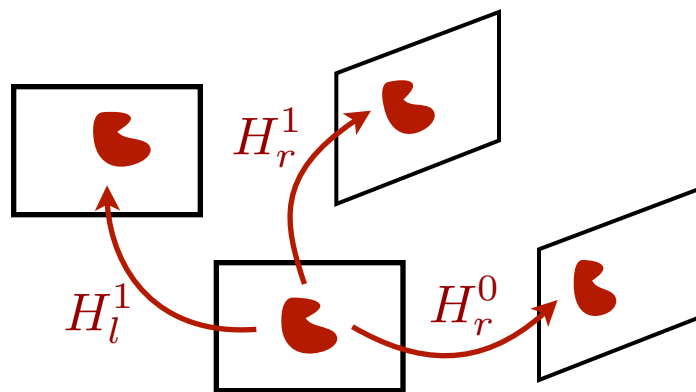
- Explicitly model scene as a collection of planar patches
- Recent trend in both
  - stereo [e.g. Bleyer et al. 2011] and
  - optical flow [e.g. Sun, Sudderth, Black 2010]



[Bleyer, Gelautz, Rother, Rhemann 2009]

# Piecewise planarity and rigidity

- For 3D scene flow, represent scene as a collection of rigidly moving planes
- Why that?
  - a good approximation for most object surfaces
  - stronger regularization
  - large, well-delineated support regions for estimation
  - simple mapping with homographies
  - (potential for implicit (or even explicit) object segmentation)



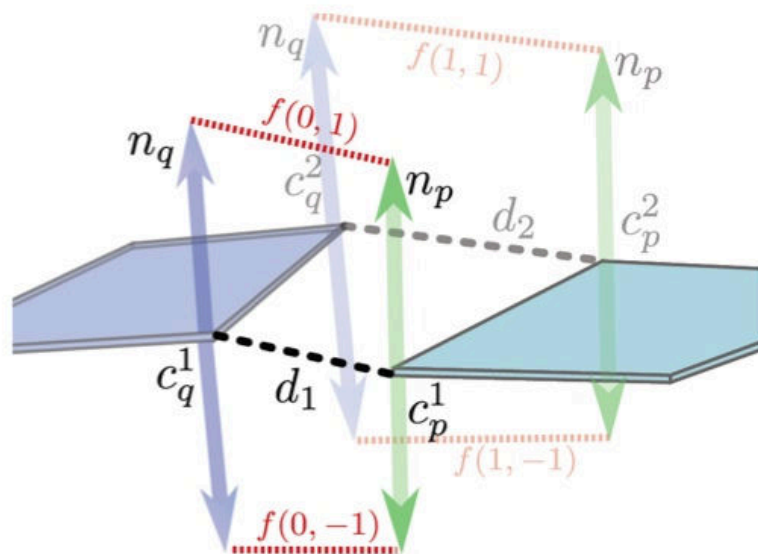
[Vogel, Roth, Schindler ICCV'13]

# Energy

- In addition to data fidelity and smoothness, encourage “good” segmentation

$$E(\mathcal{P}, \mathcal{S}) = E_D(\mathcal{P}, \mathcal{S}) + \lambda E_R(\mathcal{P}, \mathcal{S}) + \mu E_S(\mathcal{S})$$

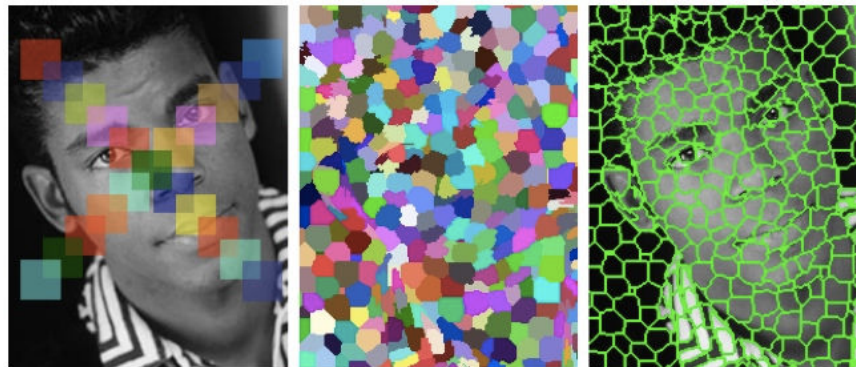
- Data term
  - two optical flow pairs
  - two stereo pairs
- Smoothness at boundary pixels
  - 3D distances between patches small
  - curvature small (computed from distances of auxiliary points)



# Segmentation regularization

- Segments should be **spatially coherent** ( $\neq$  compact)
- Suitable models exist in energy-based segmentation
  - Potts model to encourage segment boundaries at high gradients
  - allow only assignments to nearby segments

$$E_S(\mathcal{S}) = \sum_{\substack{(\mathbf{p}, \mathbf{q}) \in \mathcal{N}, \\ \mathcal{S}(\mathbf{p}) \neq \mathcal{S}(\mathbf{q})}} \exp\left(\frac{-a|I_l^0(\mathbf{p}) - I_l^0(\mathbf{q})|}{\sigma_I(\mathbf{p}, \mathbf{q}) + \epsilon}\right) + \sum_{\mathbf{p} \in I_l^0} \begin{cases} 0, & \exists \mathbf{e} \in \mathcal{E}(s_i) : \|\mathbf{e} - \mathbf{p}\|_\infty < N_S \\ \infty, & \text{else.} \end{cases}$$



[Veksler, Boykov, Mehrani 2010]



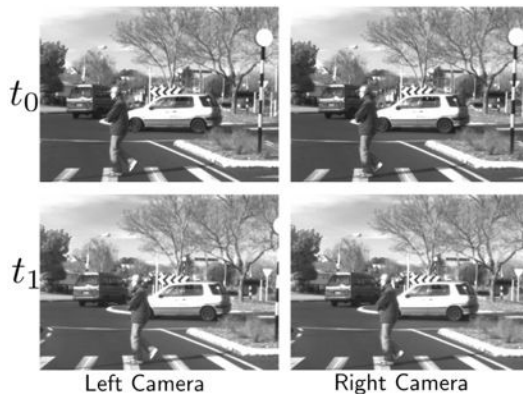
# Inference

- Initial segmentation based only on intensity
- Compute multiple scene flow proposals for each segment
  - run simpler scene flow, 2D flow + stereo
  - fit rigidly moving planes to segments



# Inference

- Assign each segment to one of the proposals
  - minimize energy function over all pixels of all segments
  - $\alpha$ -expansion with QPBO



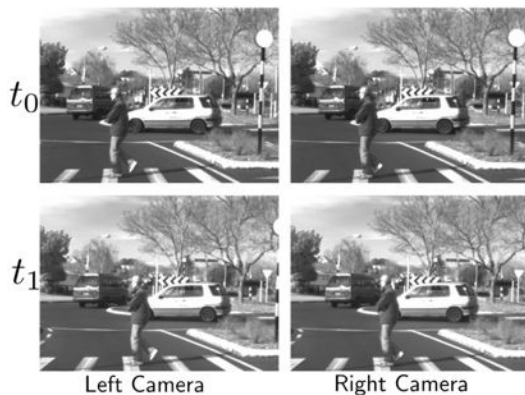
# Occlusions

- Scene flow is 3D, allows for explicit occlusion modeling
  - occluded pixels do not pay a data penalty
  - instead they are assigned a fixed occlusion penalty
- Integration into inference scheme
  - in each binary expansion step, find cases “if super-pixel  $\mathbf{p}$  is on plane  $\mathbf{X}$  and super-pixel  $\mathbf{q}$  is on plane  $\mathbf{Y}$ , then  $\mathbf{p}$  occludes  $\mathbf{q}$ ”
  - there can be  $>1$  segments on the line of sight to  $\mathbf{q}$   $\rightarrow$  higher-order cliques
  - reduce to binary cliques with auxiliary variables (carefully, to introduce few non-submodular cliques)



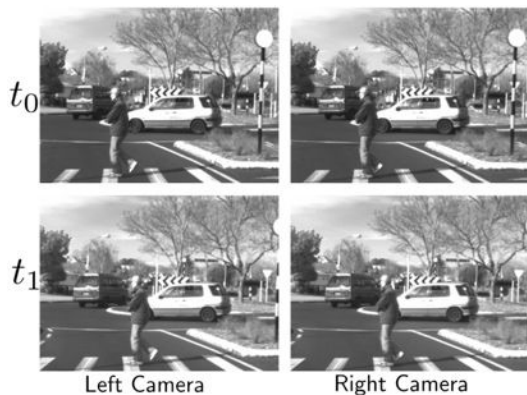
# Inference

- Estimate segment-level occlusions
  - include occlusion potentials
  - solve assignment again



# Inference

- Keep plane parameters fixed, reassign pixels to segments
  - simpler now, because only few segments within the allowed distance
- (iterate)


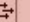


# Results

- Quantitative results
  - KITTI flow+stereo benchmark
  - errors of reprojected 2D disparity/flow (no other ground truth)

## Optical Flow Evaluation

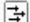

This table ranks general optical flow methods, performing a full 2D search, as com

Rank	Method	Setting	Code	Out-Noc	Out-All	Avg-Noc	Avg-All	Density
1	<a href="#">PR-Sf+E</a>	 		4.08 %	7.79 %	0.9 px	1.7 px	100.00 %

Anonymous submission

2	<a href="#">PCBP-Flow</a>	 		4.08 %	8.70 %	0.9 px	2.2 px	100.00 %
---	---------------------------	---	--	--------	--------	--------	--------	----------


K. Yamaguchi, D. McAllester and R. Urtasun: [Robust Monocular Epipolar Flow Estimation](#). CVPR 2013.

3	<a href="#">MotionSLIC</a>	 		4.36 %	10.91 %	1.0 px	2.7 px	100.00 %
---	----------------------------	---	--	--------	---------	--------	--------	----------

K. Yamaguchi, D. McAllester and R. Urtasun: [Robust Monocular Epipolar Flow Estimation](#). CVPR 2013.

4	<a href="#">PR-Sceneflow</a>	 		4.48 %	8.98 %	1.3 px	3.3 px	100.00 %
---	------------------------------	---	--	--------	--------	--------	--------	----------


Anonymous submission

5	<a href="#">TGV2ADCSIFT</a>			6.55 %	15.35 %	1.6 px	4.5 px	100.00 %
---	-----------------------------	---	--	--------	---------	--------	--------	----------

C. Vogel, S. Roth and K. Schindler: [An Evaluation of Data Costs for Optical Flow](#). German Conference on Pattern R


7	<a href="#">DeepMatching</a>			8.04 %	18.60 %	1.6 px	5.7 px	100.00 %
---	------------------------------	---	--	--------	---------	--------	--------	----------

Anonymous submission


8	<a href="#">TVL1-HOG</a>			8.31 %	19.21 %	2.0 px	6.1 px	100.00 %
---	--------------------------	---	--	--------	---------	--------	--------	----------

Anonymous submission


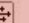
## Stereo Evaluation

Rank	Method	Setting	Code	Out-Noc	Out-All	Avg-Noc	Avg-All	Density
1	<a href="#">PCBP-SS</a>			3.49 %	4.79 %	0.8 px	1.0 px	100.00 %

K. Yamaguchi, D. McAllester and R. Urtasun: [Robust Monocular Epipolar Flow Estimation](#). CVPR 2013.

2	<a href="#">StereoSLIC</a>			3.99 %	5.17 %	0.9 px	1.0 px	99.89 %
---	----------------------------	---	--	--------	--------	--------	--------	---------


K. Yamaguchi, D. McAllester and R. Urtasun: [Robust Monocular Epipolar Flow Estimation](#). CVPR 2013.

3	<a href="#">PR-Sf+E</a>	 		4.09 %	4.95 %	0.9 px	1.0 px	100.00 %
---	-------------------------	---	--	--------	--------	--------	--------	----------

Anonymous submission

4	<a href="#">PCBP</a>			4.13 %	5.45 %	0.9 px	1.2 px	100.00 %
---	----------------------	---	--	--------	--------	--------	--------	----------

K. Yamaguchi, T. Hazan, D. McAllester and R. Urtasun: [Continuous Markov Random Fields for Robust Stereo Esti](#)

5	<a href="#">PR-Sceneflow</a>	 		4.46 %	5.32 %	1.0 px	1.1 px	100.00 %
---	------------------------------	---	--	--------	--------	--------	--------	----------

Anonymous submission

6	<a href="#">DDS</a>			4.63 %	5.44 %	1.0 px	1.1 px	100.00 %
---	---------------------	---	--	--------	--------	--------	--------	----------

Anonymous submission

7	<a href="#">wSGM</a>			5.03 %	6.24 %	1.3 px	1.6 px	97.03 %
---	----------------------	---	--	--------	--------	--------	--------	---------

T. Robert Spangenberg and R. Rojas: [Weighted Semi-Global Matching and Center-Symmetric Census Transform](#)

8	<a href="#">ATGV</a>			5.05 %	6.91 %	1.0 px	1.6 px	100.00 %
---	----------------------	---	--	--------	--------	--------	--------	----------

R. Ranftl, T. Pock and H. Bischof: [Minimizing TGV-based Variational Models with Non-Convex Data terms](#). ICSSV

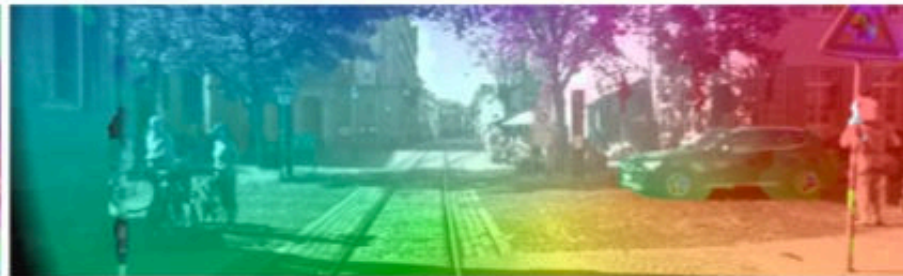
# Results

- Comparison
  - other 3D scene flow algorithms
  - intermediate stages and variants of the framework

Error threshold $Z$	FLOW ( <i>All</i> )				FLOW ( <i>Noc</i> )				STEREO( <i>All</i> )				STEREO ( <i>Noc</i> )			
	2	3	4	5	2	3	4	5	2	3	4	5	2	3	4	5
<i>LSF</i> [3]	21.6	16.9	14.3	12.7	16.0	12.0	10.0	8.8	17.6	12.0	9.0	7.2	16.4	10.8	8.0	6.3
<i>Rig</i> [24]	16.1	12.1	10.1	8.8	10.6	7.3	5.7	4.8	15.0	10.6	8.3	6.8	13.7	9.5	7.2	5.8
<i>2D</i> [11, 26]	18.9	15.0	12.8	11.3	11.0	7.9	6.5	5.7	13.5	9.9	8.0	6.7	12.3	8.9	7.0	5.8
<i>PRSSeg-3D</i>	13.8	10.1	8.2	7.1	8.4	5.6	4.5	3.9	9.4	6.8	5.4	4.6	8.4	6.0	4.8	4.0
<i>PRSPix-3D</i>	12.8	9.3	7.6	6.6	7.2	4.7	3.7	3.2	8.1	5.8	4.6	3.9	7.1	5.0	4.0	3.3
<i>PRSSeg-2D</i>	12.4	9.0	7.3	6.4	7.4	5.0	3.9	3.4	8.9	6.4	5.1	4.3	7.9	5.6	4.4	3.7
<i>PRSPix-2D</i>	11.8	8.5	6.9	6.0	6.9	4.5	3.5	3.0	8.3	5.9	4.7	3.9	7.3	5.1	4.0	3.3
<i>PRSPix-O-2D</i>	11.2	7.7	5.9	5.1	6.8	4.4	3.3	2.8	8.3	5.9	4.7	4.0	7.4	5.2	4.1	3.4
<i>PRSPix-2D+R</i>	10.9	7.6	6.0	5.1	6.3	4.1	3.1	2.7	7.9	5.7	4.5	3.8	6.9	4.8	3.8	3.2
<i>PRSPix-2D+R+E</i>	<b>10.0</b>	<b>6.7</b>	<b>5.0</b>	<b>4.1</b>	<b>5.8</b>	<b>3.6</b>	<b>2.7</b>	<b>2.2</b>	<b>7.4</b>	<b>5.3</b>	<b>4.2</b>	<b>3.5</b>	<b>6.4</b>	<b>4.5</b>	<b>3.6</b>	<b>3.0</b>

# Results

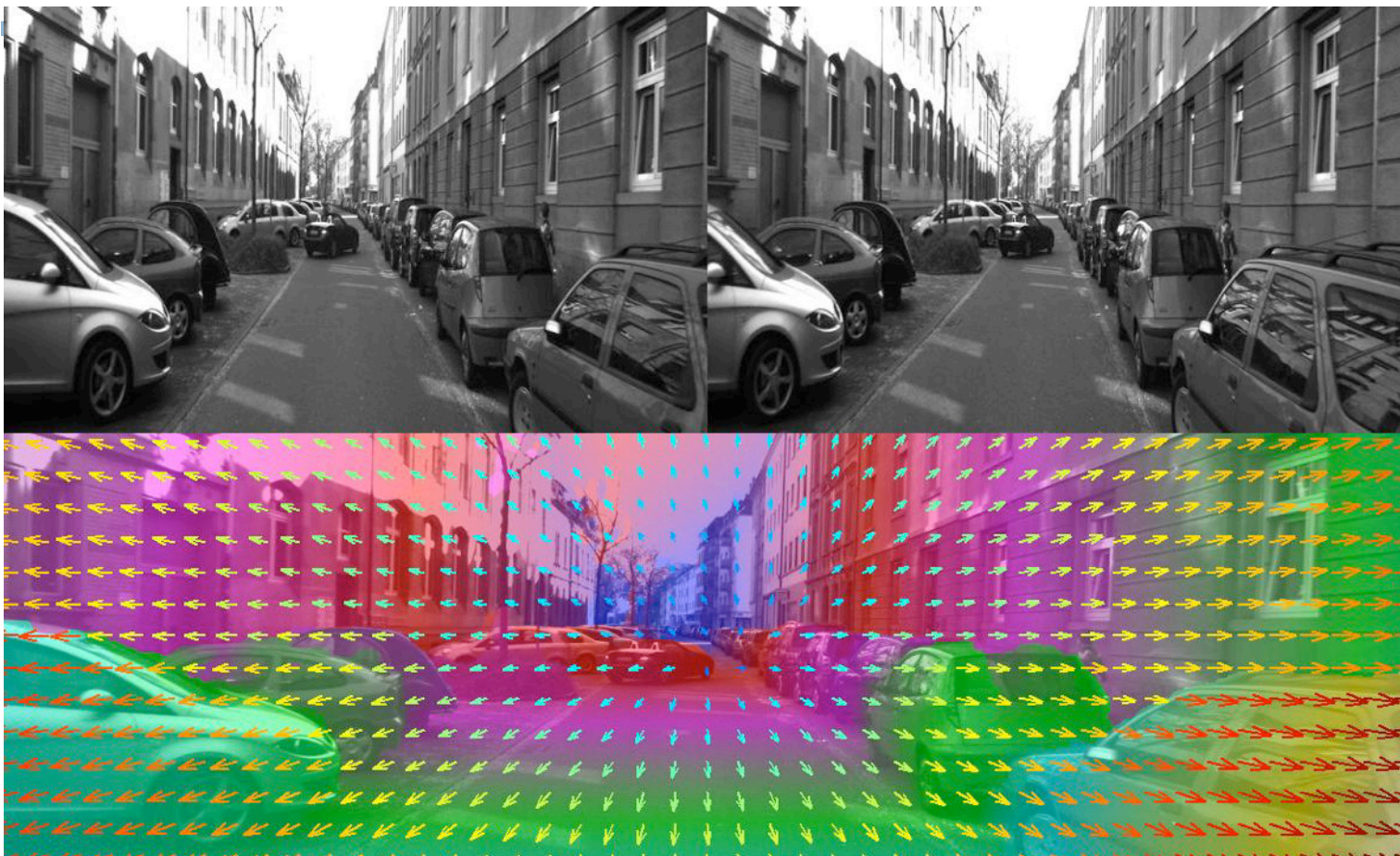
- Qualitative results





# Results

- Qualitative results



# Summary and Outlook

- Does rigidity help? **Yes!**
  - we believe it is the best regularizer for scene flow thus far
- Local or piecewise rigidity? **Hard to say;**
  - obviously depends on the scene - “horses for courses”
  - in practice piecewise so far works much better
- Something obvious missing? **Yes!**
  - >2 time steps must be beneficial (although notoriously hard to show)
  - rigidity should benefit more than agnostic smoothing
- Is rigidity the final word? **No!**
  - scene flow is still in its infancy compared to mature computer vision problems (stereo, flow, SfM, categorization...)
  - we as a community can certainly do a lot better; a good area if you want to make an impact ;-)

# Cast listing



Christoph Vogel



Stefan Roth



Konrad Schindler

# ECCV 2014 – European Conference on Computer Vision

Zurich, September 5-12, 2014

[HOME](#)

[IMPORTANT DATES](#)

[LOCAL ARRANGEMENT](#)

[PEOPLE](#)



## Important Dates

- ECCV 2014 Submission Deadline: 7 March 2014 – Friday
- ECCV 2014 Supplementary Materials Deadline: 14 March 2014 – Friday
- ECCV 2014 Announcement of Decisions : 16 June 2014 – Monday

# Photogrammetric Computer Vision - PCV 2014

ISPRS Technical Commission III Midterm Symposium

5th - 7th September 2014, Zurich, Switzerland

In Conjunction with the European Conference on Computer Vision

Home

People

Program

Submission

Important Dates

ECCV 2014



## Important dates

- submission deadline for ISPRS Annals (full paper peer-review) 13 April 2014
- submission deadline for ISPRS Archives (abstract review) 19 June 2014

isprs

information from imagery