



Motivation. Using object knowledge, we encourage disparities to agree with plausible surfaces while simultaneously recovering 3D geometry of the objects (bottom). This improves results (center) compared to current state-of-the-art stereo methods (top) [1].

#### Abstract

- Current stereo methods often fail at textureless. reflecting or semi-transparent surfaces such as cars.
- Yet, as humans we are able to effortlessly extract information about the geometry of objects even from a single image.
- Little is known about the importance of recognition for stereo matching.

We introduce object knowledge for well-constrained object categories into a slanted-plane MRF. We leverage inverse graphics to sample set of plausible object disparity maps given an initial semi-dense disparity estimate and a rough semantic segmentation of the image. We encourage the presence of these 2.5D shape samples (or "'displets") depending on how much their geometry and semantic class agrees with the image observations.

# Representation

We decompose the image into planar superpixels S. Each superpixel  $i \in S$  is associated with:

- region  $\mathcal{R}_i$  in the image
- random variable  $\mathbf{n}_i \in \mathbb{R}^3$  describing a plane in 3D

Each **displet**  $k \in \mathcal{D}$  is associated with:

• random variable  $d_k \in \{0, 1\}$  (presence/absence)

#### Our goal is to **infer**:

- superpixel planes  $\mathbf{n} = \{\mathbf{n}_i | i \in \mathcal{S}\}$
- displets  $\mathbf{d} = \{d_k | k \in \mathcal{D}\}$



**Displet Illustration.** Top: wireframe model, middle: rendered disparity map, bottom: fitted superpixel planes

We sample 3D CAD model configurations  $\xi_k$ and associate each displet  $k \in \mathcal{D}$  with

- its class label  $c_k$
- a fitness score  $\kappa_k$
- a set of superpixels  $\mathcal{S}_k$
- the corresponding planes  $\hat{\mathbf{n}}_{k,1}, \dots, \hat{\mathbf{n}}_{k,|\mathcal{S}_k|}$

#### Data term

We penalize deviations from an initial sparse disparity map  $\hat{\Omega}$  obtained using the method of [1]:

$$\varphi_i^{\mathcal{S}}(\mathbf{n}_i) = \sum_{\mathbf{p}\in\mathcal{R}_i\cap\hat{\Omega}_+} \rho_{\tau_1}(\omega(\mathbf{n}_i,\mathbf{p}) - \hat{\omega}(\mathbf{p}))$$

•  $\omega(\mathbf{n}_i, \mathbf{p})$  denotes the disparity of plane  $\mathbf{n}_i$  at pixel  $\mathbf{p}$ •  $\rho_{\tau}(\cdot)$  is a robust  $l_1$  penalty  $\rho_{\tau}(x) = \min(x, \tau)$ .

#### Local smoothness

For adjacent superpixels, we penalize discontinuities at the boundary  $\mathcal{B}_{ii}$  and encourage similar orientations:

$$\psi_{ij}^{\mathcal{S}}(\mathbf{n}_{i},\mathbf{n}_{j}) = \theta_{1} \sum_{\mathbf{p}\in\mathcal{B}_{ij}} \rho_{\tau_{2}} \left( \omega(\mathbf{n}_{i},\mathbf{p}) - \omega(\mathbf{n}_{j},\mathbf{p}) \right) + \\ \theta_{2} \rho_{\tau_{3}} \left( 1 - |\mathbf{n}_{i}^{T}\mathbf{n}_{j}| / (||\mathbf{n}_{i}|| ||\mathbf{n}_{j}||) \right)$$

# **Displets: Resolving Stereo Ambiguities using Object Knowledge**

Fatma Güney, Andreas Geiger MPI for Intelligent Systems, Tübingen

# Displets

Displets: A representative set of disparity maps for a specific semantic class (e.g., car) conditioned on the image.



**Examples of Car Displets.** Our displets cover the most likely 3D car shapes given the initial disparity map and a semantic segmentation.

# **Stereo Matching using Displets**

$$\mathbf{n}, \mathbf{d}) = \underbrace{\sum_{i \in \mathcal{S}} \varphi_i^{\mathcal{S}}(\mathbf{n}_i)}_{\text{Data}} + \underbrace{\sum_{i \sim j} \psi_{ij}^{\mathcal{S}}(\mathbf{n}_i, \mathbf{n}_j)}_{\text{Smoothness}} + \underbrace{\sum_{k \in \mathcal{D}} \varphi_k^{\mathcal{D}}(d_k)}_{\text{Displets}} + \underbrace{\sum_{k \in \mathcal{D}} \sum_{i \in \mathcal{S}_k} \psi_{ki}^{\mathcal{D}}(d_k, \mathbf{n}_i)}_{\text{Consistency}}$$

# Displet Unary

We encourage image regions with semantic class label  $c_k$ to be explained by a displet of the corresponding class:

 $\varphi_k^{\mathcal{D}}(d_k) = -\theta_3 \left[ d_k = 1 \right] \cdot \left( \left| \left[ \mathbf{S} = c_k \right] \cap \mathbf{M}_k \right| + \kappa_k \right) \right]$ 

- $\mathbf{M}_k$  is the foreground mask of displet k.
- S denotes the semantic segmentation.

# Consistency between Displets and Superpixels

We ensure consistency by encouraging superpixels to agree with the geometry of the active displets:

$$\psi_{ki}^{\mathcal{D}}(d_k, \mathbf{n}_i) = \lambda_{ki} \left[ d_k = 1 \right] \cdot \left( 1 - \delta(\mathbf{n}_i, \hat{\mathbf{n}}_{k, z_i}) \right)$$

•  $\lambda_{ki}$  denotes a weight which depends on the distance to the boundary (see experiments).

# Inference

Minimizing  $E(\mathbf{n}, \mathbf{d})$  is a non-convex, mixed continuousdiscrete optimization problem. We use max-product particle BP with TRW-S for the inner iterations.

At each outer iteration, particles are sampled

- around the previous MAP estimate
- using neighbouring superpixels

# Sampling Displet Proposals

We leverage the initial semi-dense disparity map, a semantic segmentation of the image as well as 3D CAD models to sub-sample the space of displets using MCMC.

# **Observation Model**

We sample pose parameters from a model which

- encourages displets to explain pixels in object regions
- avoids occlusion of other objects

We run one Markov chain for each combination of 3D CAD models and object proposals using Metropoliswithin-Gibbs sampling with randomly chosen blocks.

We select the 8 most dominant modes after burn-in for each model and combine all results to yield the final set of displets.

# 3D CAD Models

For sampling displets, we leverage 3D CAD models from Google Warehouse. To speed-up the rendering process, we reduce details while preserving the overall shape by fitting a "semi-convex hull" to the point clouds of the original models:



(c) QSlim

(d) MATLAB

Mesh Simplification. Comparison of our semi-convex hull algorithm to QSIim and MATLAB's reducepatch function.

# **Experimental Results on KITTI**

Illustration of Displet Influence  $\lambda_{ki}$ 

![](_page_0_Picture_72.jpeg)

**Displet Influence.** Less transparency indicates a higher penalty, e.g., we allow more deviation at the displet boundaries.

# Ablation Study

![](_page_0_Figure_75.jpeg)

#### Importance of Different Terms.

![](_page_0_Figure_77.jpeg)

Number of Proposals and Models. Influence of limiting the number of object proposals (left) and the variety of CAD models used for generating the displets (right).

### Quantitative Results

Rank	Method	Out-No	C	Out-A	411	Avg-N	loc	Avg-	-All
1	Our Method	8.40 %	%	9.89	%	1.9	рх	2.3	рх
2	VC-SF* [2]	11.58 9	%	12.29	%	ا 2.7	рх	2.8	рх
3	PCBP-SS [3]	14.26	0/0	18.33	%	2.4 j	рх	3.9	рх
4	SPS-StFI* [4]	14.74 🤅	0/0	18.00	%	2.9 j	рх	3.6	рх
I	E	1		I		I			
11	MC-CNN [1]	18.45	0/0	21.96	%	3.5 J	рх	4.3	рх
	R	eflectiv	<i>v</i> e	Regi	ons	<b>)</b> .			

Rank	Method	Out-No	c Out-All	Avg-Noc	Avg-All
1	Our Method	2.47 %	6 3.27 %	0.7 px	0.9 px
2	MC-CNN [1]	2.61 %	3.84 %	0.8 px	1.0 px
3	SPS-StFI* [4]	2.83 %	3.64 %	0.8 px	0.9 px
4	VC-SF* [2]	3.05 %	3.31 %	0.8 px	0.8 px
5	SPS-St [4]	3.39 %	4.41 %	0.9 px	1.0 px
6	PCBP-SS <sup>[3]</sup>	3.40 %	4.72 %	0.8 px	1.0 px
		' 		'	

All Regions.

![](_page_0_Picture_83.jpeg)

	Reflective	All		
	19.84 %	3.35 %		
	17.72 %	3.31 %		
	15.96 %	3.21 %		
	17.06 %	3.28 %		
	14.78 %	3.12 %		
	15.32 %	3.04 %		
	7.08 %	2.87 %		
)	7.16 %	2.78 %		

![](_page_0_Picture_87.jpeg)

References

- [1] Jure Zbontar and Yann LeCun. Computing the stereo matching cost with a convolutional neural network. In CVPR, 2015.
- [2] Christoph Vogel, Stefan Roth, and Konrad Schindler.
- View-consistent 3D scene flow estimation over multiple frames. In ECCV, 2014.
- [3] K. Yamaguchi, D. McAllester, and R. Urtasun. Robust monocular epipolar flow estimation. In CVPR, 2013.
- [4] Koichiro Yamaguchi, David McAllester, and Raquel Urtasun. Efficient joint segmentation, occlusion labeling, stereo and flow estimation. In ECCV, 2014.

# **Contact Information**

Web http://www.cvlibs.net/projects/displets Email {fatma.guney, andreas.geiger}@tue.mpg.de