

 Open access • Proceedings Article • DOI:10.1109/CVPR.2012.6247873

Layered segmentation and optical flow estimation over time — [Source link](#)

[Deqing Sun](#), [Erik B. Sudderth](#), [Michael J. Black](#)

Institutions: [Brown University](#)

Published on: 16 Jun 2012 - [Computer Vision and Pattern Recognition](#)

Topics: [Image segmentation](#), [Motion estimation](#), [Optical flow](#), [Continuous optimization](#) and [Gradient descent](#)

Related papers:

- [Large Displacement Optical Flow: Descriptor Matching in Variational Motion Estimation](#)
- [Representing moving images with layers](#)
- [A Database and Evaluation Methodology for Optical Flow](#)
- [A naturalistic open source movie for optical flow evaluation](#)
- [Motion detail preserving optical flow estimation](#)

Share this paper:    

View more about this paper here: <https://typeset.io/papers/layered-segmentation-and-optical-flow-estimation-over-time-agcqk8hd4y>

Layered Segmentation and Optical Flow Estimation Over Time: Supplemental Material

Deqing Sun¹ Erik B. Sudderth¹ Michael J. Black^{1,2}

¹Department of Computer Science, Brown University, Providence, RI 02912, USA

²Max Planck Institute for Intelligent Systems, 72076 Tübingen, Germany

Section 1 provides additional details for each of the “cooperative” moves used in optimizing the discrete layered model. Section 2 provides a high-level description of the algorithms for determining the depth ordering and the number of layers. Section 3 provides the full set of images illustrating the experimental results. The results include the screen shots of the Middlebury AAE and EPE evaluation tables at the time of writing (April 2012), as well as all experimental results on the Middlebury flow dataset and the MIT layer segmentation dataset.

1. “Cooperative” Moves for the Discrete Model

We will use a toy example to explain the effect achieved by each move. Figure 1 shows the desired layer segmentation and flow field for the input “bird apple” sequence. During optimization we will see that there are several (fairly bad) local optima and we will need to make large changes to the solution to get out of these optima. Note that the binary selection variable b controls different variables for each move. The potential functions for each move are also defined differently though the functions may share the same name.

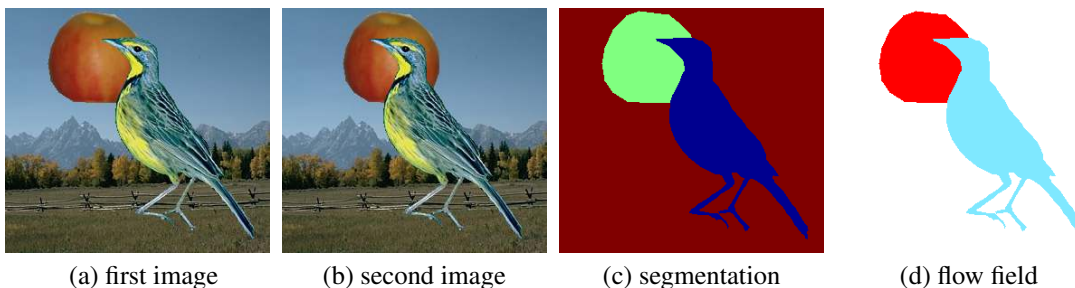


Figure 1. Input “bird apple” frames, the ground truth segmentation, and flow field.

Simultaneous segmentation and flow move. Sometimes a region may be assigned to the wrong layer but with the correct motion, as shown in Figure 2. Changing the segmentation or the flow field alone will not move the solution from the local optimum. We need to simultaneously change the segmentation and the motion. At each step, this move allows a pixel to become visible at a particular layer \hat{k} . The newly visible pixel at layer \hat{k} takes the motion of the corresponding pixel at the previously visible layer k' , while the latter takes the motion of its affine flow field. Note that the binary decision variable b_t^p may influence all the support functions for the first \hat{k} layers. When the binary decision variable is 0, the configuration is unchanged. When it is 1, $u_{t\hat{k}}^p(1) = u_{tk'}^{p,\text{old}}$, $u_{tk'}^p(1) = u_{\theta_{tk'}}^p$, $g_{tk}^p(1) = 0$, $k < \hat{k}$, and $g_{t\hat{k}}^p(1) = 1$. Using the binary variable

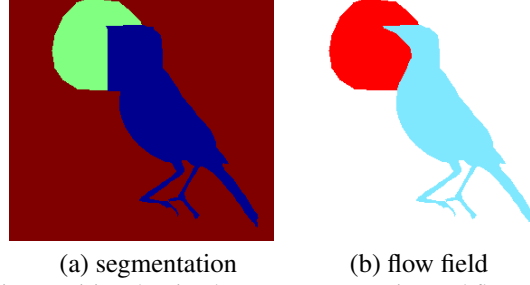


Figure 2. Local minimum of the objective requiring the simultaneous segmentation and flow move. A large region of pixels in the center has been incorrectly assigned to the bird layer though their motion is correct. Changing either the segmentation or the flow field alone will not resolve this and get us out of the local minimum. We need to make these pixels visible in the apple layer and still retain their motion using the simultaneous segmentation and flow move (see text).

representation, we can perform the following simplification

$$\begin{aligned}
E(\mathbf{b}) &= \sum_{t=1}^{T-1} \left\{ E_{\text{data}}(\mathbf{u}_t(\mathbf{b}_t), \mathbf{v}_t(\mathbf{b}_t), \mathbf{g}_t(\mathbf{b}_t), \mathbf{g}_{t+1}(\mathbf{b}_{t+1})) + \sum_{k=1}^K \lambda_a (E_{\text{space}}^{\text{flow}}(\mathbf{u}_{tk}(\mathbf{b}_t), \theta_{tk}) + E_{\text{space}}^{\text{flow}}(\mathbf{v}_{tk}(\mathbf{b}_t), \theta_{tk})) \right. \\
&\quad \left. + \sum_{k=1}^{K-1} \lambda_b E_{\text{space}}^{\text{sup}}(\mathbf{g}_{tk}(\mathbf{b}_t)) + \lambda_c E_{\text{time}}(\mathbf{g}_{tk}(\mathbf{b}_t), \mathbf{g}_{t+1,k}(\mathbf{b}_{t+1}), \mathbf{u}_{tk}(\mathbf{b}_t), \mathbf{v}_{tk}(\mathbf{b}_t)) \right\} + \sum_{k=1}^{K-1} \lambda_b E_{\text{space}}^{\text{sup}}(\mathbf{g}_{Tk}(\mathbf{b}_T)) \quad (1) \\
&= \sum_{t=1}^{T-1} \sum_p \left(\sum_{q \in \mathcal{N}_{p,\text{time}}^0} \phi_{\text{time}}^0(b_t^p, b_{t+1}^q) + \sum_{q \in \mathcal{N}_{p,\text{time}}^1} \phi_{\text{time}}^1(b_t^p, b_{t+1}^q) \right) + \sum_{t=1}^T \sum_p \left[\sum_{q \in \mathcal{N}_{p,\text{space}}} \phi_{\text{space}}(b_t^p, b_t^q) + \phi_{\text{affine}}(b_t^p) \right],
\end{aligned}$$

where $\mathbf{u}_t(\mathbf{b}_t)$ and $\mathbf{v}_t(\mathbf{b}_t)$ are deterministic functions of the candidate flow fields and the binary selection variable \mathbf{b}_t , and $\mathbf{g}_t(\mathbf{b}_t)$ is a deterministic function of the previous support function $\mathbf{g}_t^{\text{old}}$ and the binary selection variable \mathbf{b}_t .

A pixel p at each layer has two flow vector candidates to select and may interact with two different pixels at the next frame. We define different potential functions for both temporal neighborhood structures. For example, when $b_t^p = 0$, the temporal neighbors are

$$\mathcal{N}_{p,\text{time}}^0 = \{(i + [u_{tk}^p(0)], j + [v_{tk}^p(0)]), 1 \leq k \leq K\}, \quad (2)$$

and the corresponding potential function will only “fire” when $b_t^p = 0$, *i.e.*

$$\phi_{\text{time}}^0(b_t^p, b_{t+1}^q) = \left[(\rho_d(\mathbf{I}_t^p - \mathbf{I}_{t+1}^q) - \lambda_d) s_{tk}^p(b_t^p) s_{t+1,k}^q(b_{t+1}^q) + \lambda_c (1 - \delta(g_{tk}^p(b_t^p), g_{t+1,k}^q(b_{t+1}^q))) (1 - \delta(k, K)) \right] (1 - b_t^p), \quad (3)$$

which incorporates both the data term and the temporal consistency term of the first $K - 1$ support functions. We evaluate the warped image I_{t+1}^q at subpixel positions and the visibility mask $s_{t+1,k}^q$ and the warped support function $g_{t+1,k}^q$ at integer positions. Similarly

$$\mathcal{N}_{p,\text{time}}^1 = \{(i + [u_{tk}^p(1)], j + [v_{tk}^p(1)]), 1 \leq k \leq K\}, \quad (4)$$

and $\phi_{\text{time}}^1(b_t^p, b_{t+1}^q)$ only “fires” when $b_t^p = 1$

$$\phi_{\text{time}}^1(b_t^p, b_{t+1}^q) = \left[(\rho_d(\mathbf{I}_t^p - \mathbf{I}_{t+1}^q) - \lambda_d) s_{tk}^p(b_t^p) s_{t+1,k}^q(b_{t+1}^q) + \lambda_c (1 - \delta(g_{tk}^p(b_t^p), g_{t+1,k}^q(b_{t+1}^q))) (1 - \delta(k, K)) \right] b_t^p. \quad (5)$$

For the spatial pairwise term, the set $\mathcal{N}_{p,\text{space}}$ contains all the four nearest neighbors of the pixel p and is fixed throughout the document. The binary selection variable changes the states of several binary support functions and flow fields. The effects sum together as

$$\phi_{\text{space}}(b_t^p, b_t^q) = \sum_{k=1}^{K-1} \lambda_b w_q^p (1 - \delta(g_{tk}^p(b_t^p), g_{tk}^q(b_t^q))) + \sum_{k=1}^K \lambda_a \rho_{\text{mrf}}(u_{tk}^p(b_t^p) - u_{tk}^q(b_t^q)). \quad (6)$$

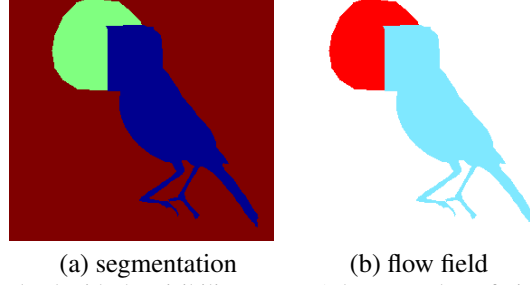


Figure 3. Local minimum that can be solved with the visibility move. A large number of pixels in the center of the image are assigned incorrectly to the bird layer, while the apple layer has the correct motion field for these pixels. Changing these pixels to be visible at the apple layer will fix this particular solution and get us out of the local minimum.

The unary term enforces the selected motion vector to be similar to the affine motion field of the newly visible layer

$$\phi_{\text{affine}}(b_t^p) = \sum_{k=1}^K \lambda_a \lambda_{\text{aff}} \rho_{\text{aff}}(u_{tk}^p(b_t^p) - u_{\theta_{tk}}^p). \quad (7)$$

We can solve this binary problem using QPBO and set the states of all these involved variables according to the binary selection variables.

Visibility move. Sometimes the segmentation is wrong, while the correct layer has the right motion field. Hence we can make big changes to the layer segmentation to correct the errors in the motion field, as shown in Figure 3. Given the current flow estimate, we decide whether to make a pixel p visible for some layer \hat{k} by modifying the previous layer support \mathbf{g}^{old} . Because the visibility state of a pixel is jointly determined by several support functions, the binary selection variable controls the states of all the support functions involved. When $b_t^p = 0$, all the support functions retain their previous value at p , *i.e.* $g_{tk}^p(0) = g_{tk}^{p,\text{old}}$. When $b_t^p = 1$, we need to adjust the support functions of the first \hat{k} layers so that layer \hat{k} is visible at p , *i.e.* $g_{tk}^p(1) = 0$ if $k < \hat{k}$, and $g_{tk}^p(1) = 1$ if $k = \hat{k}$. When \hat{k} is the last layer, all the support functions of the first $K - 1$ layers are set to be 0 at p . The energy function for the binary variable is

$$\begin{aligned} E(\mathbf{b}) &= \sum_{t=1}^{T-1} \left\{ E_{\text{data}}(\mathbf{u}_t, \mathbf{v}_t, \mathbf{g}_t(\mathbf{b}_t), \mathbf{g}_{t+1}(\mathbf{b}_{t+1})) + \lambda_c E_{\text{time}}(\mathbf{g}_{tk}(\mathbf{b}_t), \mathbf{g}_{t+1,k}(\mathbf{b}_{t+1}), \mathbf{u}_{tk}, \mathbf{v}_{tk}) \right\} + \sum_{t=1}^T \sum_{k=1}^{K-1} \lambda_b E_{\text{space}}^{\text{sup}}(\mathbf{g}_{tk}(\mathbf{b}_t)) \\ &= \sum_{t=1}^{T-1} \sum_p \sum_{q \in \mathcal{N}_{p,\text{time}}} \phi_{\text{time}}(b_t^p, b_{t+1}^q) + \sum_{t=1}^T \sum_p \sum_{q \in \mathcal{N}_{p,\text{space}}} \phi_{\text{space}}(b_t^p, b_t^q), \end{aligned} \quad (8)$$

in which $\mathcal{N}_{p,\text{time}} = \{(i + [u_{tk}^p], j + [v_{tk}^p]), 1 \leq k \leq K\}$ and the time term incorporates both the data term and the temporal consistency of the first $K - 1$ support functions

$$\phi_{\text{time}}(b_t^p, b_{t+1}^q) = \begin{cases} \left(\rho_d(\mathbf{I}_t^p - \mathbf{I}_{t+1}^{q'}) - \lambda_d \right) s_{tk}^p(b_t^p) s_{t+1,k}^q(b_{t+1}^q) + \lambda_c (1 - \delta(g_{tk}^p(b_t^p), g_{t+1,k}^q(b_{t+1}^q))), & k < K, \\ \left(\rho_d(\mathbf{I}_t^p - \mathbf{I}_{t+1}^{q'}) - \lambda_d \right) s_{tk}^p(b_t^p) s_{t+1,k}^q(b_{t+1}^q), & k = K, \end{cases} \quad (9)$$

where the visibility mask s depends on the support function which in turn depends on the binary variable b , and the corresponding pixel q depends on the flow vector of the k th layer. The visibility move may change several support functions and the potential function for the spatial term is

$$\phi_{\text{space}}(b_t^p, b_t^q) = \sum_{k=1}^{K-1} \lambda_b w_q^p \left(1 - \delta(g_{tk}^p(b_t^p), g_{tk}^q(b_t^q)) \right). \quad (10)$$

Support function move. Given the current flow estimate, we decide whether to make the support function of a pixel p at the selected layer \hat{k} equal to 1 or retain its previous value. The energy function and the potential functions are the same as Eqs. (8)-(10) but the binary variable selects the solutions in a different way. When $b_t^p = 0$, $g_{t\hat{k}}^p(0) = g_{t\hat{k}}^{p,\text{old}}$. When $b_t^p = 1$, $g_{t\hat{k}}^p(1) = 1$.

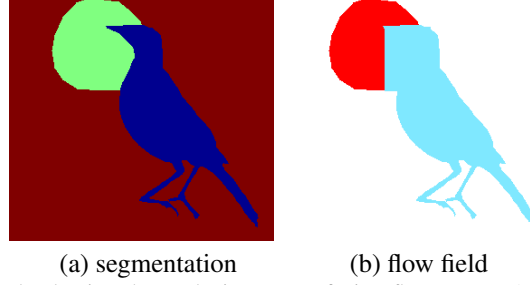


Figure 4. Local minimum that can be solved using the occlusion-aware fusion flow move. A large number of pixels in the center of the image have the wrong motion vectors. Making large changes to the flow field can reduce the errors. Note that we use the segmentation information to reason about occlusions.

Occlusion-aware fusion flow move. Here we deal with the case in which the segmentation is correct, but the motion field has a large region of errors, particularly in occlusion regions, as shown in Figure 4. We must make large changes to the flow field to reduce the motion errors. Given the current estimate of the support functions for each layer, we want to select the motion of each pixel between its current flow estimate and its affine mean flow field, $u_{tk}^p(0) = u_{tk}^{p,\text{old}}$ and $u_{tk}^p(1) = u_{\theta_{tk}}^p$. Note that for occluded pixels, their data likelihood term does not provide useful information to estimate the motion. Instead we can predict the motion of the occluded pixels using the affine motion field fitted to the visible pixels of the same object (layer). For a particular layer k of frame t ($t < T$), the energy function to minimize is

$$\begin{aligned}
 E(\mathbf{b}_{tk}) &= E_{\text{data}}(\mathbf{u}_t(\mathbf{b}_{tk}), \mathbf{v}_t(\mathbf{b}_{tk}), \mathbf{g}_t, \mathbf{g}_{t+1}) + \sum_{k=1}^K \lambda_a (E_{\text{space}}^{\text{flow}}(\mathbf{u}_{tk}(\mathbf{b}_{tk}), \theta_{tk}) + E_{\text{space}}^{\text{flow}}(\mathbf{v}_{tk}(\mathbf{b}_{tk}), \theta_{tk})) \\
 &\quad + \lambda_c \sum_{k=1}^{K-1} E_{\text{time}}(\mathbf{g}_{tk}, \mathbf{g}_{t+1,k}, \mathbf{u}_{tk}(\mathbf{b}_{tk}), \mathbf{v}_{tk}(\mathbf{b}_{tk})) \\
 &= \sum_p \phi_{\text{unary}}(b_{tk}^p) + \sum_p \sum_{q \in \mathcal{N}_{p,\text{space}}} \phi_{\text{space}}(b_{tk}^p, b_{tk}^q), \tag{11}
 \end{aligned}$$

where the unary term incorporates the data term, the temporal consistency of the support function, and the deviation from the affine motion field

$$\phi_{\text{unary}}(b_{tk}^p) = \left(\rho_d (\mathbf{I}_t^p - \mathbf{I}_{t+1}^{q'}) - \lambda_d \right) s_{tk}^p + \lambda_c (1 - \delta(g_{tk}^p, g_{t+1,k}^q)) (1 - \delta(k, K)) + \lambda_a \lambda_{\text{aff}} \rho_{\text{aff}} (u_{tk}^p(b_{tk}^p) - u_{\theta_{tk}}^p), \tag{12}$$

in which the corresponding pixel at the next frame depends on the flow field selection at the current frame, $q' = (i + u_{tk}^p(b_{tk}^p), j + v_{tk}^p(b_{tk}^p))$ and $q = (i + [u_{tk}^p(b_{tk}^p)], j + [v_{tk}^p(b_{tk}^p)])$. This energy function differs from the FusionFlow method [3] in that the segmentation information directly modulates the data likelihood term in Eq. (12) and enables occlusion reasoning.

Continuous flow refinement. Given the current estimate of the support functions for each layer, we refine the flow field for each layer by minimizing $E(\mathbf{u}_{tk}, \mathbf{v}_{tk}) =$

$$E_{\text{data}}(\mathbf{u}_{tk}, \mathbf{v}_{tk}, \mathbf{g}_t, \mathbf{g}_{t+1}) + \lambda_a (E_{\text{space}}^{\text{flow}}(\mathbf{u}_{tk}, \theta_{tk}) + E_{\text{space}}^{\text{flow}}(\mathbf{v}_{tk}, \theta_{tk})) + \lambda_c E_{\text{time}}(\mathbf{g}_{tk}, \mathbf{g}_{t+1,k}, \mathbf{u}_{tk}, \mathbf{v}_{tk}) (1 - \delta(k, K)). \tag{13}$$

Compared with standard optical flow formulation, this energy function contains the segmentation information necessary to reason about occlusions and an additional temporal consistency term for the support functions. We can still optimize the energy function using standard warping-based incremental estimation methods [4] and the refinement step adaptively changes the flow field for the discrete optimization.

2. Pseudo Code for the Discrete Optimizer

The high-level algorithm for the initialization and the discrete optimization is

```

Perform flow estimation using "Classic+NL"
Cluster initial flow into different affine motion groups

```

```
Perform Sequence A for forward and backward depth ordering
Pick the solution with lower energy
Perform Sequence B to decide depth ordering
Perform Sequence A to refine segmentation and flow
```

Sequence A tries to merge layers with similar motion together and includes

```
Perform visibility move for each layer
Remove redundant layers
Perform joint segmentation and flow move
Remove redundant layers
Perform support function move
Remove redundant layers
Re-estimate affine flow field for each layer
Perform flow selection move
Perform continuous flow refinement
```

During the optimization, some layers may have no visible pixels and are redundant in the energy minimization. Removing these layers will always reduce the energy of the solution. The new solution can explain the image data equally well but does not pay any cost for the spatial and temporal penalty terms for the redundant layers. Hence we remove these redundant layers and change the number of layers accordingly.

Sequence B determines the local depth ordering between neighboring layers and includes

```
For iter = 1:maxIters
  Select candidate pairs to compare by comparing the occlusion area
  current_solution = input_solution
  local_minimum = true
  For each selected pair of selected neighboring layers
    Swap depth ordering
    Perform visibility move for the two selected layers
    Perform support function move for the two selected layers
    If energy(new solution) < energy(current solution)
      current_solution = new_solution
      local_minimum = false
    End if
  End for
  Remove redundant layers
  If local_minimum == true
    break
  End if
End for
```

maxIters is set to be 10 in the current implementation and the algorithm usually stops after 4 – 6 iterations.

The pseudo code for selecting the candidate neighboring layer pairs is

```
Input: max_pairs (max number of layer pairs to compare)
Output: pairs of neighboring layers to compare
Compute (num_pairs) neighboring pairs using current segmentation
If num_pairs <= max_pairs
  Output all the neighboring pairs and exit
Else
  Compute occlusion regions between each neighboring pairs
  Order the pairs by the occlusion area in descending order
  Output the top max_pairs neighboring layer pairs
End of if
```

3. Additional Experimental Results

Figure 5 provides a screen shot of the top-performing methods on the Middlebury evaluation website at the time of writing (April 2012). Figure 6 shows the color keys for the ordering of layers and the flow field. Figures 7-10 show the first frame, estimated layer segmentation, and the estimated flow fields on the Middlebury optical flow training and test sets. Figure 11 shows the first frame, the segmentations by **HGVS** [2], **Layers++** [5], and **nLayers**, the flow field estimated by **nLayers**, and the human labeled ground truth.

Average endpoint error	avg. rank	Army (Hidden texture)			Mequon (Hidden texture)			Schefflera (Hidden texture)			Wooden (Hidden texture)			Grove (Synthetic)			Urban (Synthetic)			Yosemite (Synthetic)			Teddy (Stereo)				
		GT im0 im1			GT im0 im1			GT im0 im1			GT im0 im1			GT im0 im1			GT im0 im1			GT im0 im1							
		all	disc	untext	all	disc	untext	all	disc	untext	all	disc	untext	all	disc	untext	all	disc	untext	all	disc	untext	all	disc	untext		
IROF++ [63]	7.8	0.08	0.23	0.07	0.21	0.68	0.17	0.28	0.63	0.19	0.15	0.73	0.09	0.60	0.89	0.42	0.43	1.08	0.31	0.10	0.12	0.12	0.47	0.98	0.68		
Layers++ [38]	8.0	0.08	0.21	0.07	0.19	0.56	0.17	0.20	0.40	0.18	0.13	0.58	0.07	0.48	1.07	0.33	0.47	1.01	0.33	0.15	0.14	0.24	0.27	0.46	0.88	0.72	
MDP-Flow2 [40]	8.3	0.09	0.23	0.07	0.16	0.52	0.13	0.22	0.46	0.17	0.17	0.93	0.09	0.65	0.98	0.43	0.29	0.91	0.26	0.11	0.13	0.17	0.10	0.51	1.11	0.72	
nLayers [61]	8.5	0.07	0.19	0.06	0.22	0.59	0.19	0.25	0.54	0.20	0.15	0.84	0.08	0.53	0.78	0.34	0.44	1.08	0.30	0.13	0.13	0.20	0.17	0.47	0.97	0.67	
Sparse-NonSparse [59]	11.8	0.08	0.23	0.07	0.22	0.73	0.18	0.28	0.64	0.19	0.14	0.71	0.08	0.67	1.09	0.48	0.49	1.06	0.32	0.14	0.24	0.11	0.28	0.49	0.98	0.73	
COFM [64]	12.0	0.08	0.26	0.06	0.18	0.62	0.14	0.30	0.74	0.19	0.15	0.86	0.07	0.79	1.14	0.74	0.35	0.87	0.28	0.14	0.24	0.12	0.28	0.49	0.94	0.71	
Efficient-NL [66]	12.0	0.08	0.22	0.06	0.23	0.73	0.18	0.32	0.75	0.18	0.14	0.72	0.08	0.60	0.88	0.43	0.57	1.11	0.35	0.14	0.24	0.13	0.25	0.48	0.90	0.63	
TC-Flow [48]	12.5	0.07	0.21	0.06	0.15	0.59	0.11	0.31	0.87	0.20	0.14	0.16	0.86	0.75	1.11	0.54	0.42	0.84	0.25	0.11	0.12	0.29	0.43	0.62	1.35	0.93	
LSM [41]	12.7	0.08	0.23	0.07	0.22	0.73	0.18	0.28	0.64	0.19	0.14	0.70	0.09	0.66	0.97	0.48	0.50	1.06	0.33	0.15	0.15	0.12	0.29	0.50	1.11	0.99	
Ramp [68]	12.8	0.08	0.24	0.07	0.21	0.74	0.18	0.27	0.62	0.19	0.15	0.71	0.09	0.66	0.97	0.49	0.51	1.09	0.34	0.15	0.15	0.12	0.30	0.48	0.96	0.72	
Classic-NL [31]	14.3	0.08	0.23	0.07	0.22	0.74	0.18	0.29	0.65	0.19	0.15	0.73	0.09	0.64	0.93	0.47	0.52	1.12	0.33	0.16	0.16	0.13	0.29	0.49	0.98	0.74	
IROF-TV [56]	15.9	0.09	0.25	0.08	0.22	0.77	0.19	0.30	0.70	0.19	0.18	0.22	0.93	0.27	1.11	0.24	0.73	1.04	0.56	0.44	0.44	1.69	0.37	0.31	0.09	0.11	0.12

Average angle error	avg. rank	Army (Hidden texture)			Mequon (Hidden texture)			Schefflera (Hidden texture)			Wooden (Hidden texture)			Grove (Synthetic)			Urban (Synthetic)			Yosemite (Synthetic)			Teddy (Stereo)			
		GT im0 im1			GT im0 im1			GT im0 im1			GT im0 im1			GT im0 im1			GT im0 im1			GT im0 im1						
		all	disc	untext	all	disc	untext	all	disc	untext	all	disc	untext	all	disc	untext	all	disc	untext	all	disc	untext	all	disc	untext	
nLayers [61]	5.7	2.80	7.42	2.20	2.71	7.24	2.55	2.61	6.24	2.45	2.30	12.74	1.16	2.30	3.02	1.70	2.62	6.95	2.09	2.29	4.36	1.89	1.38	3.06	1.29	
Layers++ [38]	9.2	3.11	8.22	2.79	2.43	7.02	2.24	2.43	5.77	2.18	2.13	9.71	1.15	2.35	3.02	1.96	3.81	11.4	3.22	2.74	29	4.01	31	2.35	1.45	3.05
IROF++ [63]	10.1	3.17	8.69	2.61	2.79	9.61	2.33	3.43	8.86	2.38	2.87	12	1.52	2.74	3.57	2.19	3.20	9.70	2.71	1.96	7	3.45	9	1.22	1.80	4.06
MDP-Flow2 [40]	10.4	3.32	8.76	2.85	2.18	7.47	1.85	2.77	6.95	2.06	3.25	22	1.59	2.87	11	3.73	3.15	11.1	2.65	2.04	8	3.64	16	1.60	1.88	11
Sparse-NonSparse [59]	11.1	3.14	8.75	2.76	3.02	10.6	2.43	3.45	12.96	2.36	2.66	8	1.42	2.85	10	3.75	3.28	9.40	2.73	2.42	18	3.31	26	2.69	1.47	3.07
Efficient-NL [66]	11.2	3.01	8.29	2.30	3.12	10.3	2.40	3.83	9.97	2.08	2.76	10	1.45	2.64	3.51	2.07	3.06	8.23	2.49	2.53	22	3.73	17	2.46	1.91	12
LSM [41]	12.0	3.12	8.62	2.75	3.00	10.5	2.44	3.43	8.85	2.35	2.66	8	1.44	2.82	3.68	2.36	3.38	11	2.81	2.69	27	3.52	14	2.84	1.59	3.38
Ramp [68]	12.6	3.18	10.83	2.73	2.89	18	10.1	3.27	7.43	2.38	2.74	9	1.46	2.82	3.69	2.29	3.37	10	3.31	2.62	25	3.38	7	3.19	1.54	3.21
TC-Flow [48]	12.7	2.91	8.00	2.34	2.18	8.77	1.52	3.84	20	10.7	3.13	18	1.46	2.78	3.73	1.96	3.08	4	11.4	1.94	6	3.43	3.20	4.06	23	7.04
COFM [64]	14.0	3.17	8.90	2.4	2.41	8.34	1.92	3.77	17	10.5	2.71	14	1.19	3.08	15	3.92	3.83	20	10.9	2.20	12	3.35	2.91	1.62	2.56	1
Classic-NL [31]	14.3	3.20	8.72	2.81	3.02	10.6	2.44	3.46	13	8.84	2.78	11	1.43	2.83	3.68	2.31	3.40	12	9.09	2.87	34	3.82	21	2.86	1.67	3.53
SimpleFlow [52]	17.0	3.35	9.20	2.98	3.18	28	10.7	5.06	28	12.6	2.95	14	1.58	2.91	12	3.79	3.59	16	9.49	2.39	16	3.46	10	2.24	1.60	7

Figure 5. Screen shots of the Middlebury evaluation EPE/AE tables. At the time of writing (April 2012), the proposed method, **nLayers**, is ranked fourth in EPE and first in AAE. Note that the EPE ranks for the top four methods are very close; while the AAE rank of **nLayers** is significantly higher than the other methods.

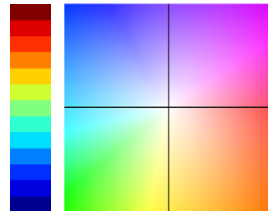
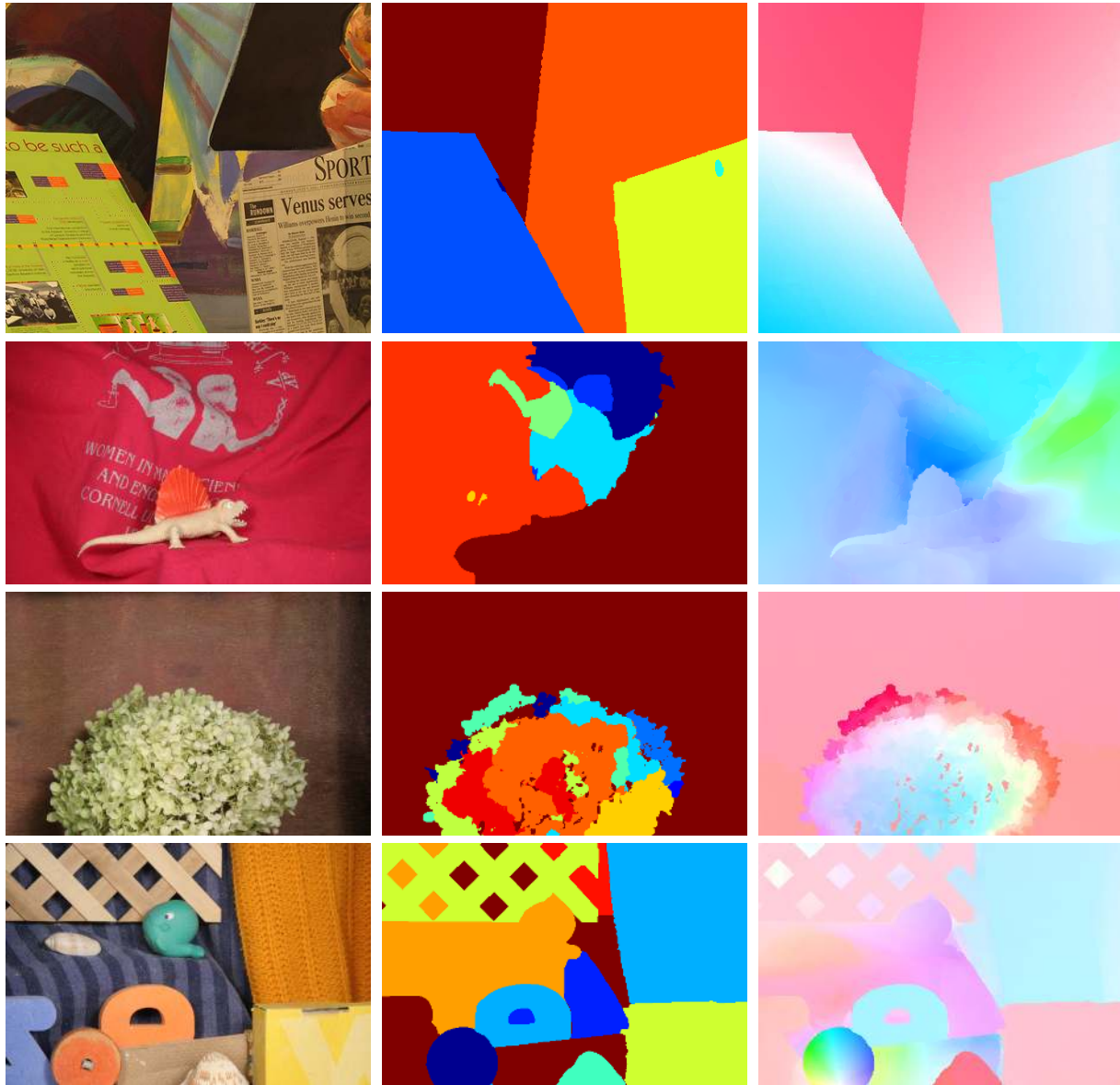


Figure 6. Left: the color key for the ordering of depth layers (blue is close and red is far); right: the color key for the flow fields [1].

References

- [1] S. Baker, D. Scharstein, J. P. Lewis, S. Roth, M. Black, and R. Szeliski. A database and evaluation methodology for optical flow. *IJCV*, 92(1):1–31, Mar. 2011. 6
- [2] M. Grundmann, V. Kwatra, M. Han, and I. Essa. Efficient hierarchical graph-based video segmentation. *CVPR*, pp. 2141–2148, 2010. 6, 11
- [3] V. Lempitsky, S. Roth, and C. Rother. FusionFlow: Discrete-continuous optimization for optical flow estimation. *CVPR*, pp. 1–8, 2008. 4
- [4] D. Sun, S. Roth, and M. Black. Secrets of optical flow estimation and their principles. *CVPR*, pp. 2432–2439, 2010. 4
- [5] D. Sun, E. Sudderth, and M. Black. Layered image motion with explicit occlusions, temporal consistency, and depth ordering. *NIPS*, pp. 2226–2234, 2010. 6, 11

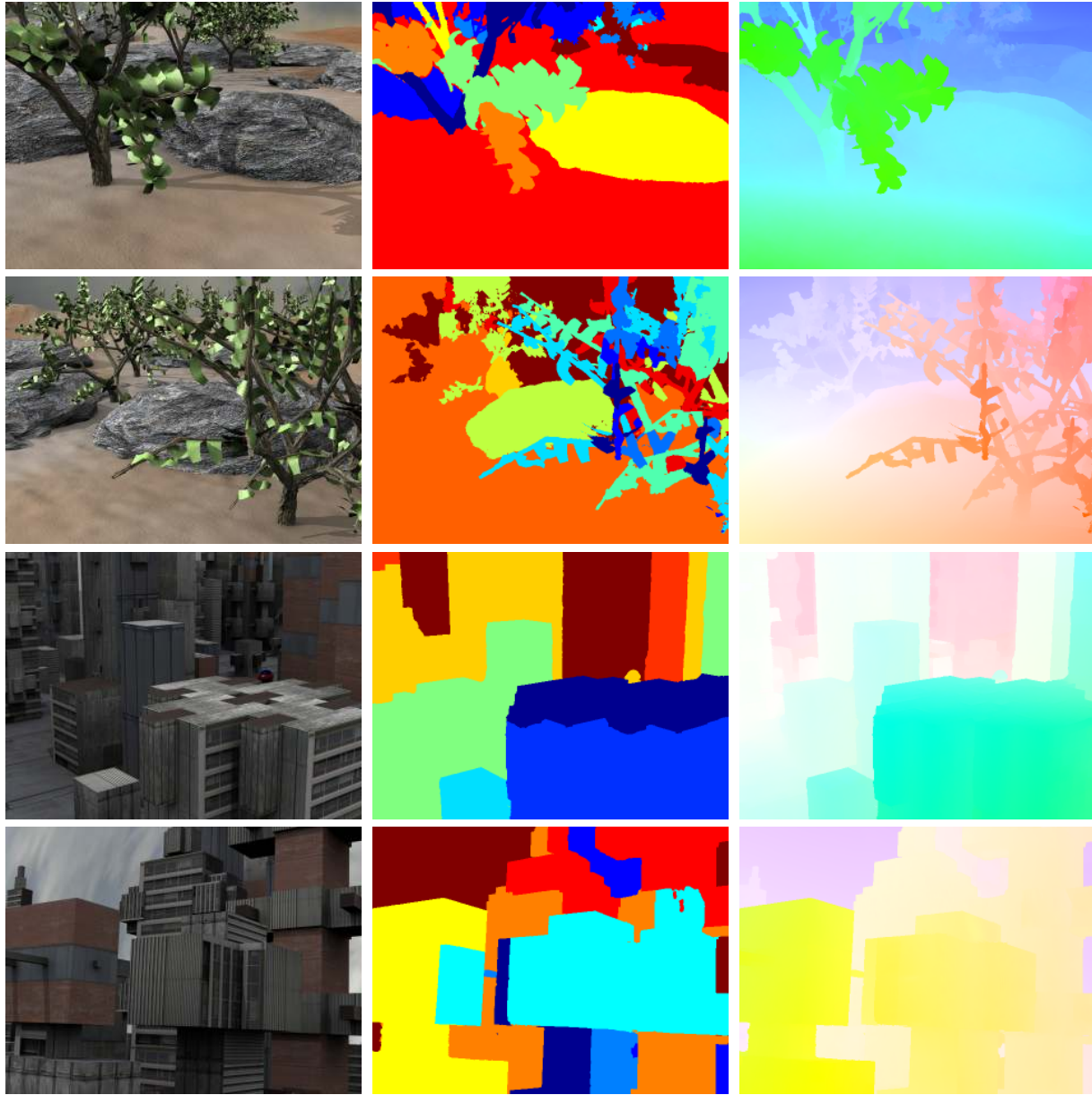


(a) First frame

(b) Layer segmentation

(c) Estimated flow field

Figure 7. Estimated flow fields and scene structure on the Middlebury *training* sequences. Left to right: first frame, layer segmentation, and estimated flow field. Top to bottom: "Venus", "Dimetrodon", "Hydrangea", and "RubberWhale".



(a) First frame

(b) Layer segmentation

(c) Estimated flow field

Figure 8. Estimated flow fields and scene structure on the Middlebury *training* sequences. Left to right: first frame, layer segmentation, and estimated flow field. Top to bottom: "Grove2", "Grove3", "Urban2", and "Urban3".

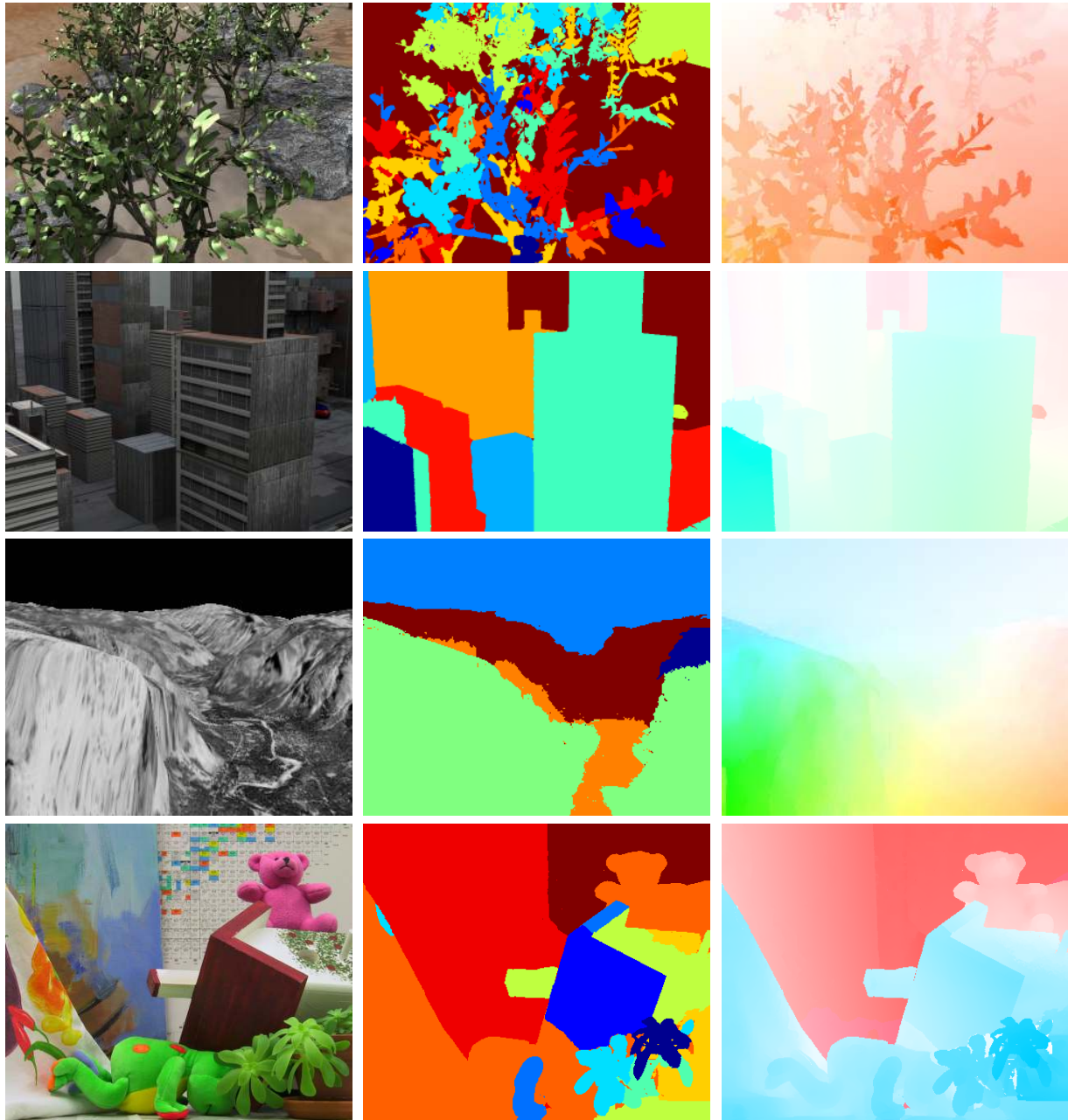


(a) First frame

(b) Layer segmentation

(c) Estimated flow field

Figure 9. Estimated flow fields and scene structure on the Middlebury *test* sequences. Left to right: first frame, layer segmentation, and estimated flow field. Top to bottom: “Army”, “Mequon”, “Schefflera”, and “Wooden”.

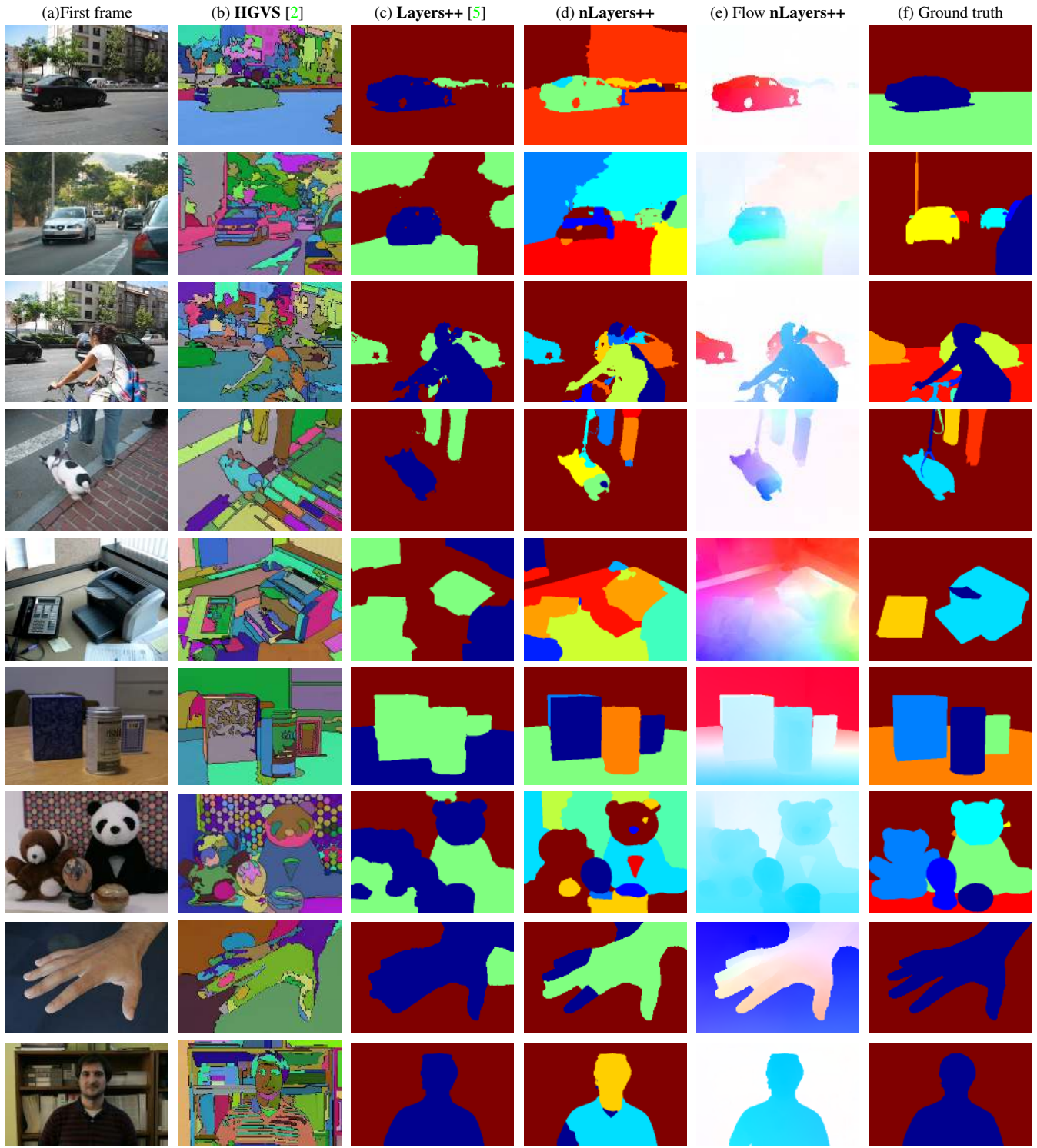


(a) First frame

(b) Layer segmentation

(c) Estimated flow field

Figure 10. Estimated flow fields and scene structure on the Middlebury *test* sequences. Left to right: first frame, layer segmentation, and estimated flow field. Top to bottom: “Grove”, “Urban”, “Yosemite”, and “Teddy”.



(a) First frame (b) HGVS [2] (c) Layers++ [5] (d) nLayers++ (e) Flow nLayers++ (f) Ground truth

Figure 11. MIT dataset. Left to right first frame, segmentation results by HGVS [2], Layers++ [5], nLayers, estimated flow field by nLayers, and human labeled ground truth. Top to bottom “car”, “car2”, “car2”, “dog”, “phone”, “table”, “toy”, “hand”, and “person”.