

Sparse Occlusion Detection with Optical Flow

Alper Ayvaci · Michalis Raptis · Stefano Soatto

Received: 12 July 2010 / Accepted: 5 August 2011
© Springer Science+Business Media, LLC 2011

Abstract We tackle the problem of detecting occluded regions in a video stream. Under assumptions of Lambertian reflection and static illumination, the task can be posed as a variational optimization problem, and its solution approximated using convex minimization. We describe efficient numerical schemes that reach the global optimum of the relaxed cost functional, for any number of independently moving objects, and any number of occlusion layers. We test the proposed algorithm on benchmark datasets, expanded to enable evaluation of occlusion detection performance, in addition to optical flow.

Keywords Occlusion detection · Optical flow · Convex optimization · Sparse optimization · Nesterov's algorithm · Split-Bregman method

1 Introduction

Occlusion phenomena are a critical component of the image formation process, shaping the statistics of natural images. Occlusion detection plays an important role in priming visual recognition of detached objects (Ayvaci and

Soatto 2011), in navigation and interaction with natural environments, and more in general in the visual information-gathering process: The complexity of the image in the “discovered (disoccluded) region” is related to the *Actionable Information Increment* in visual exploration (Soatto 2011).

Occlusions arise when a portion of the scene is visible in one image, but not another. In Da Vinci Steropsis, portions of the scene that are visible from the left eye are not visible from the right eye, and vice-versa. In a video-stream, occlusions typically occur at depth discontinuities. We are interested in determining the *occluded regions*, that is the subset of an image domain that back-projects onto portions of the scene that are not *co-visible* from a temporally adjacent image.¹ The occluded region is, in general, multiply connected, and can be quite complex, as the example of a barren tree illustrates.

Portions of the scene that are *co-visible* can be mapped onto one another by a domain deformation (Sundaramoorthi et al. 2009), called *optical flow*. It is, in general, different from the *motion field*, that is the projection onto the image plane of the spatial velocity of the scene (Verri and Poggio 1989), unless three conditions are satisfied: (a) Lambertian reflection, (b) constant illumination, and (c) constant visibility properties of the scene. Most surfaces with benign reflectance properties (diffuse/specular) can be approximated as Lambertian almost everywhere under sparse illuminants (e.g., the sun). In any case, widespread violation of Lambertian reflection does not enable correspondence, so most optical flow methods embrace (a), either implicitly or explicitly. Similarly, constant illumination (b) is a reasonable

A. Ayvaci and M. Raptis contributed equally to this work.

A. Ayvaci (✉) · M. Raptis · S. Soatto
University of California, Los Angeles, 405 Hilgard Ave, Boelter Hall #3811, Los Angeles, 90095 CA, USA
e-mail: ayvaci@cs.ucla.edu

M. Raptis
e-mail: mraptis@cs.ucla.edu

S. Soatto
e-mail: soatto@ucla.edu

¹This process could be generalized to global co-visibility, resulting in a model of the world with topologically distinct “layers” (Wang and Adelson 1994). This is beyond the scope of this paper and has already been addressed in a variational setting (Jackson et al. 2005, 2008).

assumption for ego-motion (the scene is not moving relative to the light source), and even for objects moving (slowly) relative to the light source. Assumption (c) is needed in order to have a dense flow field: If an image contains portions of its domain that are *not visible* in another image, these can patently *not* be mapped onto it by optical flow; (c) is often assumed because optical flow is defined in the limit where two images are sampled infinitesimally close in time, in which case *there are no occluded regions*, and one can focus solely on discontinuities² of the motion field. Thus, the great majority of variational motion estimation approaches provide an estimate of a dense flow field, defined at each location on the image domain, *including occluded regions*. In their defense, it can be argued that for small parallax (slow-enough motion, or far-enough objects, or fast-enough temporal sampling) occluded areas are small. However, small does not mean absent, nor unimportant, as occlusions are critical to perception (Gibson 1984) and a key for developing representations for recognition. For this reason, *we focus on occlusion detection* in video streams.

Occlusion detection would be easy if the motion field was known. Vice-versa, optical flow estimation would be easy if the occluded domain was known. As often in vision problems, one knows neither, so in the process of inferring the object of inference (the occluded domain) we will estimate optical flow, the “nuisance variable,” as a byproduct.

In this manuscript we (I) show that, starting from the standard assumptions (a)–(b), the problem of detecting (multiply-connected) occlusion regions can be formulated as a variational optimization problem (Sect. 2). We then (II) show how the functional to be minimized can be relaxed into a sequence of convex functionals and minimized using *re-weighted* ℓ_1 optimization (Appendix A and (16)). At each iteration, the functional to be minimized is related to those used for optical flow estimation, but the minimization is with respect to the indicator function of the occluded region, not just the (dense) optical flow field. We then bring to bear two different approaches to optimize these functionals, one is (III) an optimal first-order method, due to Nesterov (Sect. 3), and the other is (IV) an alternating minimization technique, known as split-Bregman method (Sect. 4). We evaluate our approach empirically in Sects. 5 and 1.2, and discuss its strengths and limitations of in Sect. 6.

A preliminary conference version of this paper has appeared in Ayvaci et al. (2010). The current version provides an additional optimization method, split-Bregman, for improving the computation speed. We have also included

the experiments analyzing re-weighting steps and compared proposed method to robust flow estimation methods.

1.1 Prior Related Work

Several algorithms have been proposed to perform occlusion detection. Many define occlusions as the regions where forward- and backward-motion are inconsistent (Proesmans et al. 1994; Alvarez et al. 2007). This is problematic as the motion in the occluded region is not just inconsistent, it is *undefined*, as there is no “motion” (domain deformation) that takes one image onto another. Other approaches (Lim et al. 2002; Kolmogorov and Zabih 2001; Sun et al. 2005) formulate occlusion directly as a classification problem, and perform motion estimation in a discrete setting, where it is an NP-hard problem. This can then be approximated with combinatorial optimization. Others (Ince and Konrad 2008; Ben-Ari and Sochen 2007) also exploit motion symmetry to detect occlusions and weight the inconsistencies with a monotonically decreasing function.

The residual from optical flow estimation has also been used to decide whether a region is occluded. Strecha et al. (2004) proposed a probabilistic formulation to detect occluded regions using the estimated noise model and histogram of occluded pixel intensities. Xiao et al. (2006) threshold the residual obtained using level-set methods to find occluded areas. Both try to minimize a non-convex energy function by iterating between two subproblems, occlusion detection and motion estimation.

Another set of algorithms infer occlusion boundaries (Stein and Hebert 2009; He and Yuille 2010) and occluded regions (Humayun et al. 2011) training a learning based detector using appearance, motion and depth features. The accuracy of these methods largely depends on the performance of the underlying feature detectors.

Occlusions have also been a concern in the optical flow community since the first global formulation was proposed by Horn and Schunck (1981). Black and Anandan (1996) proposed replacing the ℓ_2 norm of the residual with a non-convex Lorentzian penalty. Another common criterion used for this purpose is the ℓ_1 norm of the residual, which is non-trivial to minimize since it is non-smooth. Bruhn et al. (2005) and Brox et al. (2004) use Charbonnier’s penalty that is a differentiable approximation of the ℓ_1 norm; others (Wedel et al. 2008, 2009, Werlberger et al. 2009) solved the non-smooth problem with primal-dual methods decoupling the matching and regularization terms. However, none of these robust flow estimation methods focus on the detection of occlusions.

1.2 Evaluation

Optical flow estimation is a mature area of computer vision, and benchmark datasets have been developed, the best

²In occluded regions, the problem is *not* that optical flow is discontinuous; it is simply not defined; *it does not exist*. Motion in occluded regions can be *hallucinated* or *extrapolated*, based on the prior or regularizer. However, whatever motion is assigned to an occluded region *cannot be validated from the data*.

known example being the Middlebury (Baker et al. 2007). Unfortunately, the Middlebury benchmark does not provide ground truth for occlusions on the entire dataset or a scoring mechanism to evaluate the performance of occlusion detection algorithms. Unfortunately, this also biases the motion estimation scoring mechanism as ground truth motion is provided on the entire image domain, *including occluded regions*, where it can be *extrapolated* using the priors/regularizers, but not validated from the data.

To overcome this gap, we have produced a new benchmark by taking a subset of the training data in the Middlebury dataset, and hand-labeling occluded regions. We then use the same evaluation method of the Middlebury *for the (ground truth) regions that are co-visible* in at least two images. This provides a motion estimation score. Then, we provide a separate score for occlusion detection, in terms of precision-recall curves. This dataset (that at the moment is limited by our ability to annotate occluded regions to a subset of the full Middlebury, but that we will continue to expand over time), as well as the implementation of our algorithm in source format has been released publicly (<http://vision.ucla.edu/optical-flow>).

2 Joint Occlusion Detection and Optical Flow Estimation

In this section and in Appendix A, we show how the assumptions (a)–(b) can be used to formulate occlusion detection and optical flow estimation as a joint optimization problem. We assemble a functional that penalizes the (unknown) optical flow residual in the (un-known) co-visible regions, as well as the area of the occluded region. The resulting optimization problem has to be solved jointly with respect to the unknown optical flow field, and the indicator function of the occluded region.

Let $I : D \subset \mathbb{R}^2 \times \mathbb{R}^+ \rightarrow \mathbb{R}^+$; $(x, t) \mapsto I(x, t)$ be a grayscale time-varying image defined on a domain D . Under the assumptions (a)–(b), the relation between two consecutive frames in a video $\{I(x, t)\}_{t=0}^T$ is given by

$$I(x, t) = \begin{cases} I(w(x, t), t + dt) + n(x, t), & x \in D \setminus \Omega(t; dt), \\ \rho(x, t), & x \in \Omega(t; dt), \end{cases} \quad (1)$$

where $w : D \times \mathbb{R}^+ \rightarrow \mathbb{R}^2$; $x \mapsto w(x, t) \doteq x + v(x, t)$ is the domain deformation mapping $I(x, t)$ onto $I(x, t + dt)$ everywhere *except at occluded regions*. Usually *optical flow* denotes the incremental displacement $v(x, t) \doteq w(x, t) - x$. The occluded region Ω can change over time depending on the temporal sampling interval dt and is not necessarily simply-connected; so even if we call Ω the occluded region (singular), it is understood that it can be made of several disconnected components. Inside Ω , the image can take

any value $\rho : \Omega \times \mathbb{R}^+ \rightarrow \mathbb{R}^+$ that is in general unrelated to $I(w(x, t + dt))|_{x \in \Omega}$. In the limit $dt \rightarrow 0$, $\Omega(t; dt) = \emptyset$. Because of (almost-everywhere) continuity of the scene and its motion (i), and because the additive term $n(x, t)$ compounds the effects of a large number of independent phenomena³ and therefore we can invoke the Law of Large Numbers (ii), in general we have that

$$(i) \quad \lim_{dt \rightarrow 0} \Omega(t; dt) = \emptyset, \quad \text{and} \quad (ii) \quad n \stackrel{IID}{\sim} \mathcal{N}(0, \lambda), \quad (2)$$

i.e., the additive uncertainty is normally distributed in space and time with an isotropic and small variance $\lambda > 0$. We define the residual $e : D \rightarrow \mathbb{R}$ on the entire image domain $x \in D$, via

$$e(x, t; dt) \doteq I(x, t) - I(w(x, t), t + dt) = \begin{cases} n(x, t), & x \in D \setminus \Omega, \\ \rho(x, t) - I(w(x, t), t + dt), & x \in \Omega, \end{cases} \quad (3)$$

which we can write as the sum of two terms, $e_1 : D \rightarrow \mathbb{R}$ and $e_2 : D \rightarrow \mathbb{R}$, also defined on the entire domain D in such a way that

$$\begin{cases} e_1(x, t; dt) \doteq \rho(x, t) - I(w(x, t), t + dt), & x \in \Omega, \\ e_2(x, t; dt) \doteq n(x, t), & x \in D \setminus \Omega. \end{cases} \quad (4)$$

Note that e_2 is undefined in Ω , and e_1 is undefined in $D \setminus \Omega$, in the sense that they can take any value there, including zero, which we will assume henceforth. We can then write, for any $x \in D$,

$$I(x, t) = I(w(x, t), t + dt) + e_1(x, t; dt) + e_2(x, t; dt) \quad (5)$$

and note that, because of (i) e_1 is *large but sparse*, while because of (ii) e_2 is *small but dense*.⁴ We will use this as an inference criterion for w , seeking to optimize a data fidelity term that minimizes the number of nonzero elements of e_1 (a proxy of the area of Ω), and the negative log-likelihood of n .

$$\begin{aligned} \psi_{\text{data}}(w, e_1) &\doteq \|e_1\|_{\mathbb{L}^0(D)} + \frac{1}{\lambda} \|e_2\|_{\mathbb{L}^2(D)} \quad \text{subject to (5)} \\ &= \frac{1}{\lambda} \|I(x, t) - I(w(x, t), t + dt) \\ &\quad - e_1\|_{\mathbb{L}^2(D)} + \|e_1\|_{\mathbb{L}^0(D)}, \end{aligned} \quad (6)$$

³ $n(x, t)$ collects all unmodeled phenomena including deviations from Lambertian reflection, illumination changes, quantization error, sensor noise, and later also linearization error. It does *not* capture occlusions, since those are explicitly modeled.

⁴*Sparse* stands for almost everywhere zero on D . Similarly, *dense* stands for almost everywhere non-zero.

where $\|f\|_{\mathbb{L}^0(D)} \doteq |\{x \in D | f(x) \neq 0\}|$ and $\|f\|_{\mathbb{L}^2(D)} \doteq \int_D |f(x)|^2 dx$. Unfortunately, we do not know anything about e_1 other than the fact that it is sparse, and that what we are looking for is $\chi(\Omega) \propto e_1$, where $\chi : D \rightarrow \mathbb{R}^+$ is the characteristic function that is non-zero when $x \in \Omega$, i.e., where the occlusion residual is non-zero. So, the data fidelity term depends on w but also on the characteristic function of the occlusion domain Ω . For a sufficiently small dt , we can approximate⁵ for any $x \in D \setminus \Omega$,

$$I(x, t + dt) = I(x, t) + \nabla I(x, t)v(x, t) + n(x, t), \quad (9)$$

where the linearization error has been incorporated into the uncertainty term $n(x, t)$. Therefore, following the same previous steps, we have

$$\psi_{\text{data}}(v, e_1) = \|\nabla I v + I_t - e_1\|_{\mathbb{L}^2(D)} + \lambda \|e_1\|_{\mathbb{L}^0(D)}. \quad (10)$$

Since we typically do not know the variance λ of the process n , we will treat it as a tuning parameter, and because ψ_{data} or $\lambda \psi_{\text{data}}$ yield the same minimizer, we have attributed the multiplier λ to the second term. In addition to the data term, because the unknown v is infinite-dimensional and the problem is ill-posed, we need to impose regularization, for instance by requiring that the total variation (TV) be small

$$\psi_{\text{reg}}(v) = \mu \|v_1\|_{TV} + \mu \|v_2\|_{TV} \quad (11)$$

where v_1 and v_2 are the first and second components of the optical flow v , μ is a multiplier factor to weight the strength of the regularizer and the weighted isotropic TV norm is defined by

$$\|f\|_{TV(D)} = \int_D \sqrt{(g_x(x) \nabla_x f(x))^2 + (g_y(x) \nabla_y f(x))^2} dx,$$

where g_x and g_y are given by

$$g_x(x) = \exp(-\zeta \|\nabla_x I(x)\|_2) + \nu, \quad (12)$$

$$g_y(x) = \exp(-\zeta \|\nabla_y I(x)\|_2) + \nu, \quad (13)$$

where ν is small constant, preventing g_x and g_y to take the value 0 and ζ is a normalizing factor. TV is desirable in

⁵In a digital image, both domains D and Ω are discretized into a lattice, and dt is fixed. Therefore, spatial and temporal derivative operators are approximated, typically, by first-order differences. We use the formal notation

$$\nabla I(x, t) \doteq \begin{bmatrix} I(x + \begin{bmatrix} 1 \\ 0 \end{bmatrix}, t) - I(x, t) \\ I(x + \begin{bmatrix} 0 \\ 1 \end{bmatrix}, t) - I(x, t) \end{bmatrix}^T, \quad (7)$$

$$I_t(x, t) \doteq I(x, t + dt) - I(x, t). \quad (8)$$

the context of occlusion detection because it does not penalize motion discontinuities significantly. The overall problem can then be written as the minimization of the cost functional $\psi = \psi_{\text{data}} + \psi_{\text{reg}}$, which is

$$\begin{aligned} \hat{v}_1, \hat{v}_2, \hat{e}_1 = \arg \min_{v_1, v_2, e_1} & \|\nabla I v + I_t - e_1\|_{\mathbb{L}^2(D)}^2 \\ & + \lambda \|e_1\|_{\mathbb{L}^0(D)} \\ & + \mu \|v_1\|_{TV(D)} + \mu \|v_2\|_{TV(D)}. \end{aligned} \quad (14)$$

In a digital image, the domain D is quantized into an $M \times N$ lattice Λ , so we can write (14) in matrix form as:

$$\begin{aligned} \hat{v}_1, \hat{v}_2, \hat{e}_1 = \arg \min_{v_1, v_2, e_1} & \frac{1}{2} \|A[v_1, v_2, e_1]^T + b\|_{\ell_2}^2 \\ & + \lambda \|e_1\|_{\ell_0} + \mu \|v_1\|_{TV} + \mu \|v_2\|_{TV}. \end{aligned} \quad (15)$$

where $e_1 \in \mathbb{R}^{MN}$ is the vector obtained from stacking the values of $e_1(x, t)$ on the lattice Λ on top of one another (column-wise), and similarly with the vector field components $\{v_1(x, t)\}_{x \in \Lambda}$ and $\{v_2(x, t)\}_{x \in \Lambda}$ stacked into MN -dimensional vectors $v_1, v_2 \in \mathbb{R}^{MN}$. The spatial derivative matrix A is given by

$$A = [\text{diag}(\nabla_x I) \text{diag}(\nabla_y I) - \mathcal{I}],$$

where \mathcal{I} is the $MN \times MN$ identity matrix, and the temporal derivative values $\{I_t(x, t)\}_{x \in \Lambda}$ are stacked into b . For finite-dimensional vectors $u \in \mathbb{R}^{MN}$, $\|u\|_{\ell^2} = \sqrt{\langle u, u \rangle}$, $\|u\|_{\ell_0} = |\{u_i | u_i \neq 0\}|$ and $\|u\|_{TV}$ is defined as

$$\|u\|_{TV} = \sum \sqrt{((g_x)_i(u_{i+1} - u_i))^2 + ((g_y)_i(u_{i+M} - u_i))^2},$$

where g_x and g_y are the stacked versions of $\{g_x(x)\}_{x \in \Lambda}$ and $\{g_y(x)\}_{x \in \Lambda}$.

The problem (15) is NP-hard when solved with respect to the variable e_1 whose non-zero elements indicates the occluded region at each pixel in the image. A relaxation into a convex would simply replace the ℓ_0 norm with ℓ_1 . Unfortunately, this implies that “bright” occluded regions are penalized more than “dim” ones, which is clearly not desirable. Therefore, we relax the ℓ_0 norm with the *weighted- ℓ_1* norm such that

$$\begin{aligned} \hat{v}_1, \hat{v}_2, \hat{e}_1 = \arg \min_{v_1, v_2, e_1} & \frac{1}{2} \|A[v_1, v_2, e_1]^T + b\|_{\ell_2}^2 \\ & + \lambda \|W e_1\|_{\ell_1} + \mu \|v_1\|_{TV} + \mu \|v_2\|_{TV}, \end{aligned} \quad (16)$$

where W is a diagonal matrix and resort to an iterative procedure called *reweighted- ℓ_1* , proposed by Candes et al. (2008) to adapt the weights so as to better approximate the ℓ_0 norm. W is initially set to be the identity matrix, and

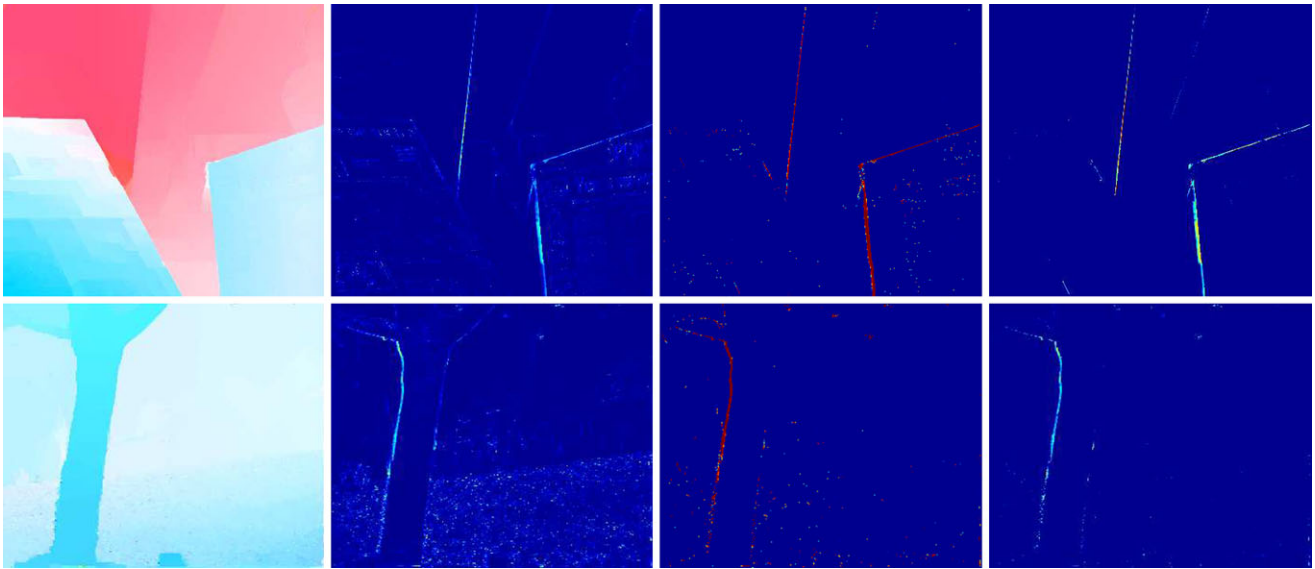


Fig. 1 The result of the proposed approach on “Venus” from Baker et al. (2007) and “Flower Garden.” The *first column* shows the motion estimates, color-coded as in Baker et al. (2007), the *second* is the

residual $I(x, t) - I(w(x), t + dt)$ before re-weighting stage; the *third* shows $|We_1|$ after re-weighting, and the *fourth* is the sparse error term e_1

correspondingly (16) is the customary convex relaxation⁶ of the original NP-hard problem (Tibshirani 1996). Each iteration has a globally optimal solution that can be reached efficiently from any initial condition. An improved approximation of the ℓ_0 norm can be obtained by adapting the weight W iteratively, for instance choosing W to be a diagonal with elements $w(x) \approx 1/(|e_1(x)| + \epsilon)$ as proposed in Candes et al. (2008). The resulting solution of (16) greatly improves sparsity, and the residual e_1 is closer to a piecewise constant (indicator) function, as shown in Fig. 1.

Note that the residual e_1 in (5) is sometimes referred to as modeling *illumination changes* (Shulman and Herve 1989; Negahdaripour 1998; Teng et al. 2005; Kim et al. 2005). However, even though the model (5) appears similar, the *priors* on e_1 are rather different. They favor smooth illumination changes; we favor sparse occlusions. While sparsity follows directly from the assumption (i), illumination changes would require a *reflectance function* to be modeled. Instead, all models of the form (5) lump reflectance and illumination into a single irradiance term (Soatto et al. 2003).

⁶This norm has been previously used in optical flow estimation, and it makes sense in that context where occlusions are the “nuisance factors.” In our context, however, occlusions are the object of inference, and we do not wish to suppress them in order to provide an optical flow reading in the occluded region, where it is undefined. Instead, we prefer using the reweighted approach as a better approximation of the original problem (16) that does not penalize bright occlusions.

3 Minimization with Nesterov’s Algorithm

In this section, we describe an efficient algorithm to solve (16) based on Nesterov’s first order scheme (Nesterov 1983) which provides $O(1/k^2)$ convergence in k iterations, whereas for standard gradient descent, it is $O(1/k)$, a considerable advantage for a large scale problem such as (16). To simplify the notation we let $(e_1)_i \doteq w_i(e_1)_i$, so that $A \doteq [\text{diag}(\nabla_x I) \text{diag}(\nabla_y I) - W^{-1}]$. The main steps of the algorithm are shown in the following table

Initialize v_1^0, v_2^0, e_1^0 . For $k \geq 0$

1. Compute $\nabla\psi(v_1^k, v_2^k, e_1^k)$
2. Compute α_k and τ_k

$$\alpha_k = 1/2(k+1), \quad \tau_k = 2/(k+3)$$

3. Compute y_k

$$y_k = [v_1^k, v_2^k, e_1^k]^T - (1/L)\nabla\psi(v_1^k, v_2^k, e_1^k)$$

4. Compute z_k

$$z_k = [v_1^0, v_2^0, e_1^0]^T - (1/L) \sum_{i=0}^k \alpha_i \nabla\psi(v_1^i, v_2^i, e_1^i)$$

5. Update $[v_1^k, v_2^k, e_1^k]^T = \tau_k z_k + (1 - \tau_k) y_k$.

Stop when the solution converges.

In order to implement this scheme, we need to address the nonsmooth nature of ℓ_1 . This is done in Nesterov (2005), that has already been used profitably for noise

reduction, inpainting and deblurring (Weiss et al. 2009; Dahl et al. 2009), and incorporated in software libraries for sparse recovery (Becker et al. 2009). In our case, we write $\psi(v_1, v_2, e_1)$ as a summation of terms

$$\begin{aligned}\psi(v_1, v_2, e_1) &= \psi_1(v_1, v_2, e_1) + \lambda \psi_2(e_1) \\ &\quad + \mu \psi_3(v_1) + \mu \psi_4(v_2)\end{aligned}$$

and compute the gradient of each term separately: The first is straightforward

$$\nabla_{v_1, v_2, e_1} \psi_1(v_1, v_2, e_1) = A^T A[v_1, v_2, e_1]^T + A^T b.$$

The other three, however, require smoothing. $\psi_2(e_1) = \|e_1\|_{\ell_1}$ can be rewritten in terms of its conjugate $\psi_2(e_1) = \max_{\|u\|_{\infty} \leq 1} \langle u, e_1 \rangle$. The smooth approximation proposed in Nesterov (2005) is

$$\psi_2^\sigma(e_1) = \max_{\|u\|_{\infty} \leq 1} \langle u, e_1 \rangle - \frac{1}{2} \sigma \|u\|_{\ell_2}^2 \quad (17)$$

which is differentiable; its gradient is u^σ , the optimal solution of (17). Consequently, $\nabla_{e_1} \psi_2^\sigma(e_1)$ is given by

$$u_i^\sigma = \begin{cases} \sigma^{-1}(e_1)_i, & |(e_1)_i| < \sigma, \\ \text{sgn}((e_1)_i), & \text{otherwise.} \end{cases} \quad (18)$$

Following the lines of Becker et al. (2009), $\nabla_{v_1} \psi_3$ is given by

$$\nabla_{v_1} \psi_3^\sigma(v_1) = G^T u^\sigma \quad (19)$$

where $G = [G_x, G_y]^T$, G_x and G_y are weighted horizontal and vertical differentiation operators, and u^σ has the form $[u_x, u_y]$ where

$$(u_{x,y})_i = \begin{cases} \sigma^{-1}(G_{x,y}v_1)_i, \\ \|[(G_x v_1)_i \ (G_y v_1)_i]^T\|_{\ell_2} < \sigma, \\ \|[(G_x v_1)_i \ (G_y v_1)_i]^T\|_{\ell_2}^{-1} (G_{x,y}v_1)_i \\ \text{otherwise.} \end{cases} \quad (20)$$

$\nabla_{v_2} \psi_4$ can also be computed in the same way. We now have all the terms necessary to compute

$$\nabla \psi(v_1, v_2, e_1) = \nabla \psi_1 + [\lambda \nabla_{e_1} \psi_2, \mu \nabla_{v_1} \psi_3, \mu \nabla_{v_2} \psi_4]^T. \quad (21)$$

We also need the Lipschitz constant L to compute the auxiliary variables y_k and z_k to minimize ψ . Since $\|G^T G\|_2$ is bounded above (Dahl et al. 2009) by 8, given the coefficients λ and μ , L is given by

$$L = \max(\lambda, 8\mu)/\sigma + \|A^T A\|_2.$$

A crucial element of the scheme is the selection of σ . It trades off accuracy and speed of convergence. A large σ

yields a smooth solution, which is undesirable when minimizing the ℓ_1 norm. A small σ causes slow convergence. We have chosen σ empirically, although the continuation algorithm proposed in Becker et al. (2009) could be employed to adapt σ during convergence.

4 Minimization with Split-Bregman

We also describe an alternative method for solving the optimization problem (16) based on the split-Bregman method proposed by Goldstein and Osher (2009), by decoupling the differentiable and non-differentiable portions of the cost function (16).

We replace $G_x v_{(1,2)}$ by $d_x^{(1,2)}$ and $G_y v_{(1,2)}$ by $d_y^{(1,2)}$ yielding to a constrained problem,

$$\begin{aligned}\hat{v}_1, \hat{v}_2, \hat{e}_1 &= \arg \min_{v_1, v_2, e_1} \frac{1}{2\mu} \|A[v_1, v_2, e_1]^T + I_t\|_{\ell_2}^2 \\ &\quad + \frac{\lambda}{\mu} \|W e_1\|_{\ell_1} \\ &\quad + \|(d_x^1, d_y^1)\|_{\ell_1} + \|(d_x^2, d_y^2)\|_{\ell_1}\end{aligned} \quad (22)$$

subject to:

$$\begin{aligned}d_x^1 &= G_x v_1, & d_y^1 &= G_y v_1, \\ d_x^2 &= G_x v_2, & d_y^2 &= G_y v_2.\end{aligned}$$

where for finite dimensional vectors $u_1, u_2 \in \mathbb{R}^{MN}$, $\|(u_1, u_2)\|_{\ell_1} = \sum_{i=1}^{MN} \sqrt{(u_1)_i^2 + (u_2)_i^2}$. By relaxing the hard constraints, the cost function (23) takes a form that can be minimized by split-Bregman:

$$\begin{aligned}\hat{v}_1, \hat{v}_2, \hat{e}_1, \hat{d}_{(x,y)}^{(1,2)} &= \arg \min_{v_1, v_2, e_1, d_{(x,y)}^{(1,2)}} \frac{1}{2\mu} \|A[v_1, v_2, e_1]^T + I_t\|_{\ell_2}^2 \\ &\quad + \frac{\lambda}{\mu} \|W e_1\|_{\ell_1} \\ &\quad + \|(d_x^1, d_y^1)\|_{\ell_1} + \|(d_x^2, d_y^2)\|_{\ell_1} \\ &\quad + \frac{\beta}{2} \|d_x^1 - G_x v_1 - b_x^1\|_{\ell_2}^2 \\ &\quad + \frac{\beta}{2} \|d_y^1 - G_y v_1 - b_y^1\|_{\ell_2}^2 \\ &\quad + \frac{\beta}{2} \|d_x^2 - G_x v_2 - b_x^2\|_{\ell_2}^2 \\ &\quad + \frac{\beta}{2} \|d_y^2 - G_y v_2 - b_y^2\|_{\ell_2}^2\end{aligned} \quad (23)$$

where β indicates the amount of relaxation. To solve (23), we divide the optimization problem into three subproblems and solve them iteratively. The first subproblem is

$$\begin{aligned} \hat{v}_1^{k+1}, \hat{v}_2^{k+1} = \arg \min_{v_1, v_2} & \frac{1}{2\mu} \|A[v_1, v_2, e_1^k]^T + I_t\|_{\ell_2}^2 \\ & + \frac{\beta}{2} \|(d_x^1)^k - G_x v_1 - (b_x^1)^k\|_{\ell_2}^2 \\ & + \frac{\beta}{2} \|(d_y^1)^k - G_y v_1 - (b_y^1)^k\|_{\ell_2}^2 \\ & + \frac{\beta}{2} \|(d_x^2)^k - G_x v_2 - (b_x^2)^k\|_{\ell_2}^2 \\ & + \frac{\beta}{2} \|(d_y^2)^k - G_y v_2 - (b_y^2)^k\|_{\ell_2}^2. \end{aligned} \quad (24)$$

The solution of this problem is straightforward. From the optimality conditions, we reach to the following system of equations

$$\begin{aligned} & \left(\frac{1}{\mu} I_x^T I_x + \beta(G_x^T G_x + G_y^T G_y) \right) v_1 + \frac{1}{\mu} I_x^T I_y v_2 \\ & = -\frac{1}{\mu} I_x^T (-e_1 + I_t) \\ & \quad + \beta \left(G_x^T (d_x^1 - b_x^1) + G_y^T (d_y^1 - b_y^1) \right) \\ & \left(\frac{1}{\mu} I_y^T I_y + \beta(G_x^T G_x + G_y^T G_y) \right) v_2 + \frac{1}{\mu} I_y^T I_x v_1 \\ & = -\frac{1}{\mu} I_y^T (-e_1 + I_t) \\ & \quad + \beta \left(G_x^T (d_x^2 - b_x^2) + G_y^T (d_y^2 - b_y^2) \right), \end{aligned}$$

where to simplify the notation we have defined the diagonal matrices $I_x = \text{diag}(\nabla_x I)$ and $I_y = \text{diag}(\nabla_y I)$. Following Goldstein and Osher (2009), to achieve efficiency, we solve the system of equations using Gauss-Seidel's method. The component-wise Gauss-Seidel solution to this problem is given by

$$\begin{aligned} (v_1)_i &= \frac{-\mu(k_2)_i(I_y)_i(I_x)_i + (k_1)_i(\mu(I_y)_i^2 + (k_3)_i)}{(k_3)_i(\mu(I_x)_i^2 + \mu(I_y)_i^2 + (k_3)_i)} \\ &= (\gamma_1)_i \\ (v_2)_i &= \frac{-\mu(k_1)_i(I_y)_i(I_x)_i + (k_2)_i(\mu(I_x)_i^2 + (k_3)_i)}{(k_3)_i(\mu(I_x)_i^2 + \mu(I_y)_i^2 + (k_3)_i)} \\ &= (\gamma_2)_i \end{aligned} \quad (25)$$

where k_1, k_2 and k_3 are given by

$$\begin{aligned} (k_1)_i &= -\mu(I_x)_i(-e_1)_i + (I_t)_i \\ & \quad + \beta \left((g_x)_{i-1}^2 (v_1)_{i-1} + (g_x)_i^2 (v_1)_{i+1} \right. \\ & \quad \left. + (g_y)_{i-M}^2 (v_1)_{i-M} + (g_y)_i^2 (v_1)_{i+M} \right) \\ & \quad + \beta \left((G_x^T d_x^1)_i + (G_y^T d_y^1)_i - (G_x^T b_x^1)_i - (G_y^T b_y^1)_i \right), \\ (k_2)_i &= -\mu(I_y)_i(-e_1)_i + (I_t)_i \\ & \quad + \beta \left((g_x)_{i-1}^2 (v_2)_{i-1} + (g_x)_i^2 (v_2)_{i+1} \right. \\ & \quad \left. + (g_y)_{i-M}^2 (v_2)_{i-M} + (g_y)_i^2 (v_2)_{i+M} \right) \\ & \quad + \beta \left((G_x^T d_x^2)_i + (G_y^T d_y^2)_i - (G_x^T b_x^2)_i - (G_y^T b_y^2)_i \right), \\ k_3 &= \beta \text{diag} (G_x^T G_x + G_y^T G_y). \end{aligned}$$

Subsequently, we need to solve the second subproblem which is given by

$$\begin{aligned} & (\hat{d}_x^{(1,2)})^{k+1}, (\hat{d}_y^{(1,2)})^{k+1} \\ & = \arg \min_{d_x^{(1,2)}, d_y^{(1,2)}} \|d_x^{(1,2)}, d_y^{(1,2)}\|_{\ell_1} \\ & \quad + \frac{\beta}{2} \|(d_x^{(1,2)}) - G_x v_{(1,2)} - (b_x^{(1,2)})^k\|_{\ell_2}^2 \\ & \quad + \frac{\beta}{2} \|(d_y^{(1,2)}) - G_y v_{(1,2)} - (b_y^{(1,2)})^k\|_{\ell_2}^2. \end{aligned} \quad (26)$$

This problem can be solved analytically using the generalized shrinkage formula (Wang et al. 2007) such that

$$\begin{aligned} (\hat{d}_x^{(1,2)})^{k+1} &= \max(s^k - 1/\beta, 0) \frac{G_x v_{(1,2)}^k + (b_x^{(1,2)})^k}{s^k}, \\ (\hat{d}_y^{(1,2)})^{k+1} &= \max(s^k - 1/\beta, 0) \frac{G_y v_{(1,2)}^k + (b_y^{(1,2)})^k}{s^k}, \end{aligned} \quad (27)$$

where $s_{(1,2)}^k$ is given by

$$(s^k)_i = \sqrt{|G_x v_{(1,2)}^k + (b_x^{(1,2)})^k|^2 + |G_y v_{(1,2)}^k + (b_y^{(1,2)})^k|^2}.$$

The remaining subproblem is

$$\hat{e}_1 = \arg \min_{e_1} \frac{1}{2\mu} \|A[v_1, v_2, e_1]^T + I_t\|_{\ell_2}^2 + \frac{\lambda}{\mu} \|W e_1\|_{\ell_1} \quad (28)$$

and can also be solved using the shrinkage operator. The solution is

$$(e_1)_i^{k+1} = \frac{r_i^k}{|r_i^k|} \max(|r_i^k| - \lambda w_i, 0), \quad (29)$$

where $r^k = I_x v_1^k + I_y v_2^k + I_t$.

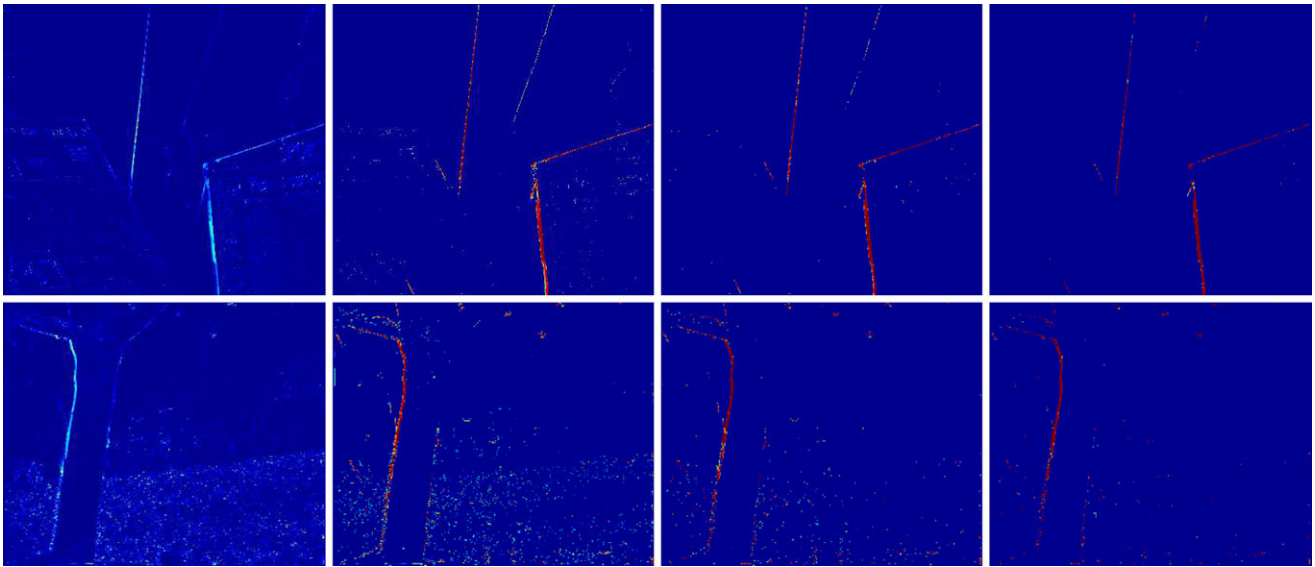


Fig. 2 This figure illustrates the initial estimate of the error term e_1 (first column) and how sparsity of $|We_1|$ improves with each of the three re-weighting iterations

The main steps of the algorithm can be summarized as follows

Initialize $v_1, v_2, e_1, d_x^1, d_y^1, d_x^2$ and d_y^2 with 0. For $k \geq 0$

$$v_1^{k+1} = \gamma_1^k, \quad v_2^{k+1} = \gamma_2^k$$

$$(d_x^{(1,2)})^{k+1} = \max(s^k - 1/\beta, 0) \frac{G_x v_{(1,2)}^k + (b_x^{(1,2)})^k}{s^k}$$

$$(d_y^{(1,2)})^{k+1} = \max(s^k - 1/\beta, 0) \frac{G_y v_{(1,2)}^k + (b_y^{(1,2)})^k}{s^k}$$

$$(b_x^{(1,2)})^{k+1} = (b_x^{(1,2)})^k + (G_x v_{(1,2)}^{k+1} - (d_x^{(1,2)})^{k+1})$$

$$(b_y^{(1,2)})^{k+1} = (b_y^{(1,2)})^k + (G_y v_{(1,2)}^{k+1} - (d_y^{(1,2)})^{k+1})$$

$$(e_1)_i^{k+1} = \frac{r_i^k}{|r_i^k|} \max(|r_i^k| - \lambda w_i, 0)$$

Stop when the solution converges.

5 Experiments

Following Sect. 1.2, evaluation of our algorithm on standard datasets is not straightforward, because these typically do not provide ground-truth occlusions. The Middlebury benchmark (Baker et al. 2007) provides occluded regions only at some parts of the dataset, so we used it as a starting point, and generated occlusion maps as follows: for each training sequence, we computed the residual given the ground truth motion and marked the regions where the resid-

ual is high. Next, we annotated the regions where ground truth is not defined. Finally, we manually fixed obvious errors in the occlusion maps. In this section, we evaluate the motion estimation and occlusion detection performance of our approach on this dataset and on the well-known Flower Garden sequence. We have also compared our algorithm to Wedel et al. (2008), Black and Anandan (1996) and Kolmogorov and Zabih (2001) quantitatively.

To handle the large motion, we run our method on a Gaussian pyramid with a scale factor 0.5 up to 5 levels. We also apply 5 warping steps at each pyramid level. In all the experiments, the coefficient λ is fixed at 0.01 while μ is increased gradually from 0.00008 to 0.01 with each warping step at each pyramid level. Relying less on the prior of the flow field at the early warping steps results in more accurate flow estimates. For the re-weighting step, we have also fixed the coefficient ϵ to 0.001. In our experiments, we also use a non-linear pre-filtering of the images to reduce the influence of illumination changes (Rudin et al. 1992; Wedel et al. 2008; Sun et al. 2010) to initialize the re-weighting stage with an accurate flow field. However, during the re-weighting steps we use the original images since pre-filtering reduces the occlusion detection accuracy.

We start with unit weights $W = \mathcal{I}$ and solve the convex problem (16) (referred as Huber- ℓ_1 in our experiments). We then adapt the weights iteratively, thus improving sparsity and achieving a better approximation of the indicator function e_1 of the occluded domain, Fig. 1. One can also observe a gradual improvement of the sparsity of $|We|$ after each re-weighting iteration, Fig. 2. At each step, the accuracy of occlusion detection also improves.

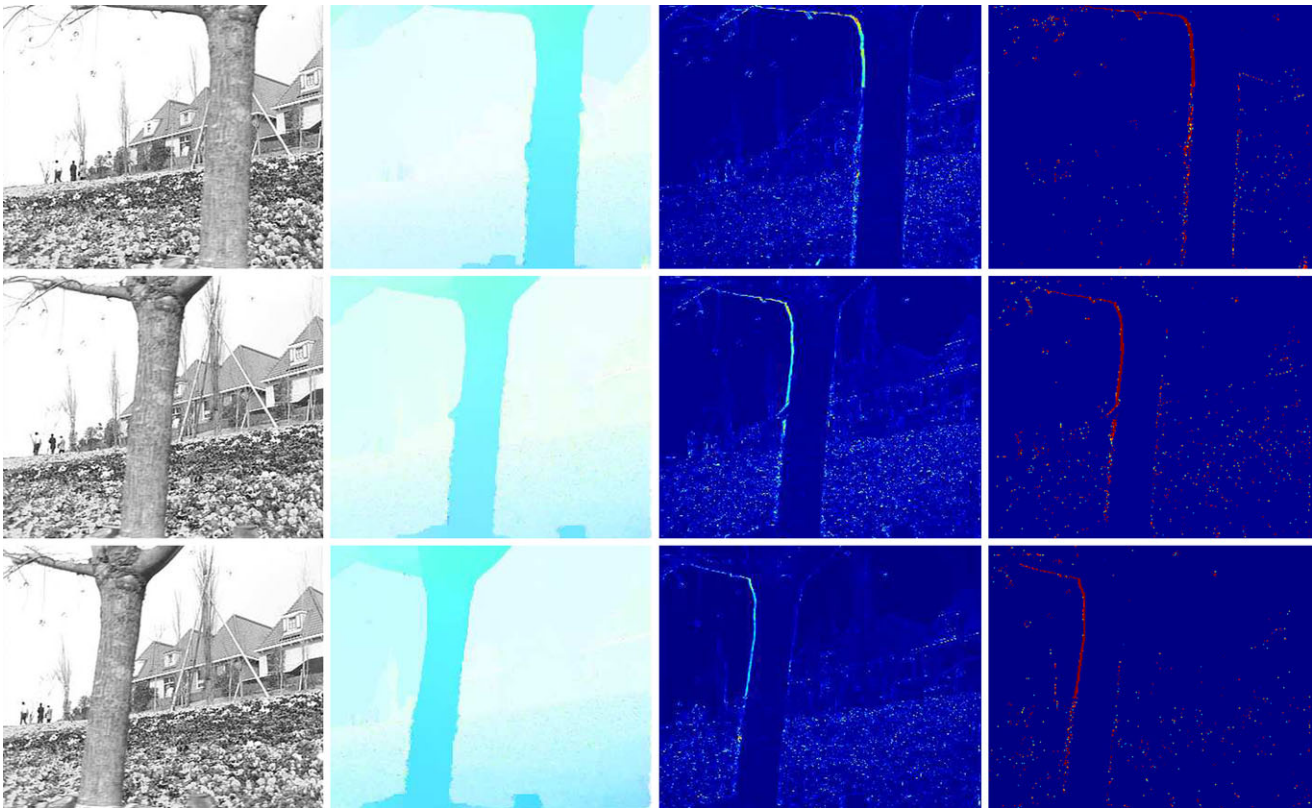


Fig. 3 Occlusion and motion estimates for more frames of the Flower Garden. Left to right: initial frame, flow estimate (*left*), initial estimate of the error term e_1 (*middle*), and occluded region (*right*)

Representative results for the Flower Garden sequence are shown in Fig. 3, where the complex occlusions produced by the foliage are also detected successfully.

In Figs. 5 and 6 we show the effects of re-weighting on the Middlebury data set. The weighted e_1 is not only sparser compared to the residual $|I(x, t) - I(w(x), t + dt)|$ computed before the re-weighting steps but also has a superior occlusion detection accuracy unlike the residual which contains regions that are not occluded. One might think that the residual could just be thresholded, instead of iteratively re-weighted. To evaluate that, we have generated precision-recall curves and observed the change of occlusion detection performance in terms of F-measure by thresholding both signals while varying the threshold value in the interval $[0, 1]$, Fig. 4. In most cases, the accuracy of the re-weighting approach is superior and more stable under the varying threshold values since $|We_1|$ better approximates an indicator function. Therefore, one can just choose non-zero elements of e_1 ⁷ to detect occluded regions instead of searching for a global threshold. Note that here the weight matrix W is the one computed at the previous re-weighting step. We have also observed that the precision-recall curves

for $|We_1|$ does not span the whole recall range, since recall value 1 is not reachable unless all the zero-elements added to the decision which is not meaningful for the analysis of a sparse signal (Fig. 4, PR-curves).

We have also compared our approach to the robust flow estimation methods proposed by Black and Anandan (Black and Anandan 1996) (using the improved version (Classical-L)) by Sun et al. (2010) and Wedel et al. (2008) by evaluating the occlusion detection accuracy on the residual $|I(x, t) - I(w(x), t + dt)|$ computed using their flow fields, Fig. 4. Furthermore, since violations of the symmetry between backward and forward flow estimations have been used intensively as an indicator of occlusions (Alvarez et al. 2007; Ince and Konrad 2008; Kolmogorov and Zabih 2001; Sun et al. 2005), we have also implemented another baseline algorithm checking such mismatch between flow fields estimated with Wedel et al. (2008). Note that neither residual nor flow field symmetry violations provides a sparse solution; therefore thresholding is required after normalization in this case. Hence, we have measured the accuracy of these techniques in terms of precision-recall and F-measure under varying threshold, Fig. 4. Our method outperforms both approaches.

⁷Notice that $w(x) > 0, \forall x$. Therefore, $We_1 \neq 0 \iff e_1 \neq 0$.

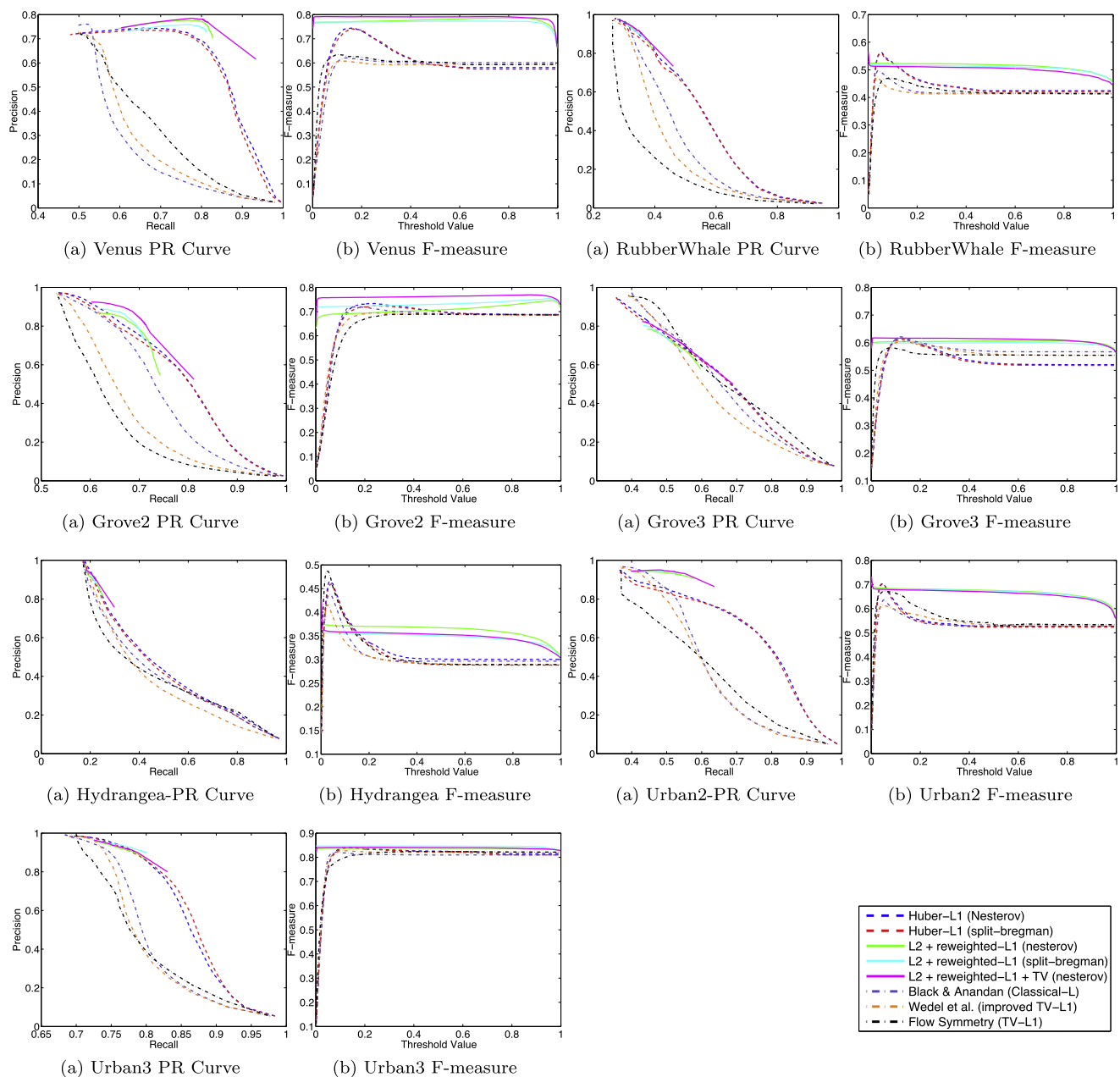


Fig. 4 Comparison of the occlusion detection accuracy of different variants of proposed technique and two other baseline methods in terms of precision-recall curves and F-measure. The first baseline method considers the residual of the optical flow estimated via Black and

Anandan (1996) or Wedel et al. (2008) as an indicator of occlusions while the second one identifies the regions of mismatch in forward and backward flow fields estimated with Wedel et al. (2008) as occluded

One might be tempted to regularize the geometry of the occluded region, for instance by adding a regularizing term $\|W_{\ell_1}\|_{TV}$ to (16). We have also evaluated this model, Fig. 4. However, occlusions can manifest themselves with very complex geometry and topology, as the Hydrangea in Figs. 4 and 6 illustrate. In such cases, a geometric regularizer is counter-productive as it generates a large number of missed detections.

We have compared our occlusion detection results to (Kolmogorov and Zabih 2001), using the code provided on-line by the authors (Table 1). Comparing motion estimates gives an unfair advantage to our algorithm because their approach is based on quantized disparity values, so the accuracy of our motion estimates is predictably superior.

We have also compared the accuracy of the solution of Nesterov's algorithm and split-Bregman's method and their

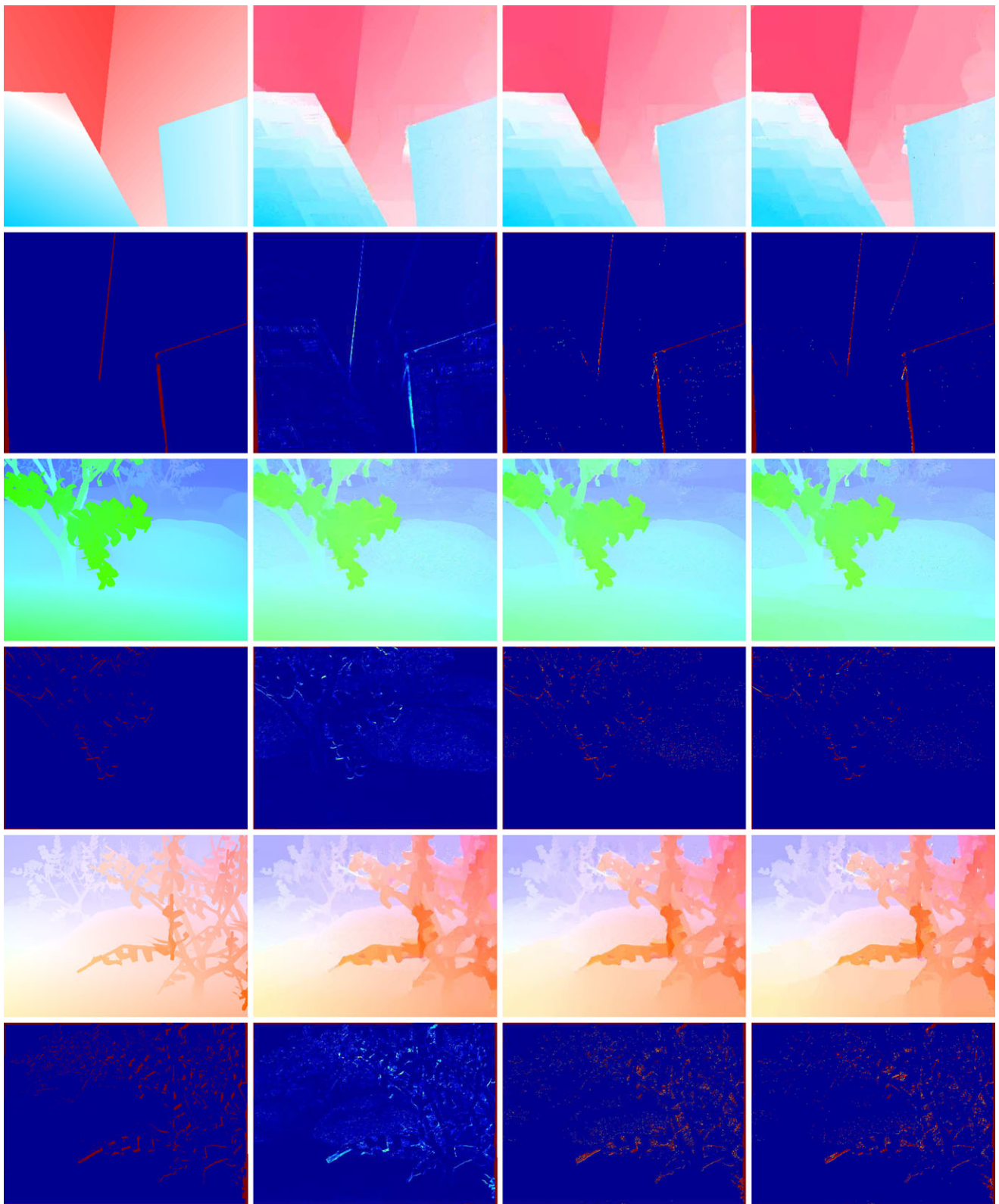


Fig. 5 This figure presents the occlusion and motion estimates on the sequences Venus, Grove2 and Grove3 from Middlebury dataset. Each sequence occupies two rows. The *odd rows*, left to right: ground truth optical flow, flow estimates before re-weighting stage and flow estimates after reweighting with Nesterov's algorithm and split-Bregman

method. The *even rows*, left to right: ground truth occluded regions, the initial estimate of the error term e_1 , the estimate of $|We_1|$ after the reweighting step with Nesterov's algorithm and split-Bregman method. The parts that are carried outside of the image domain under \hat{v} are also marked on the occlusion maps

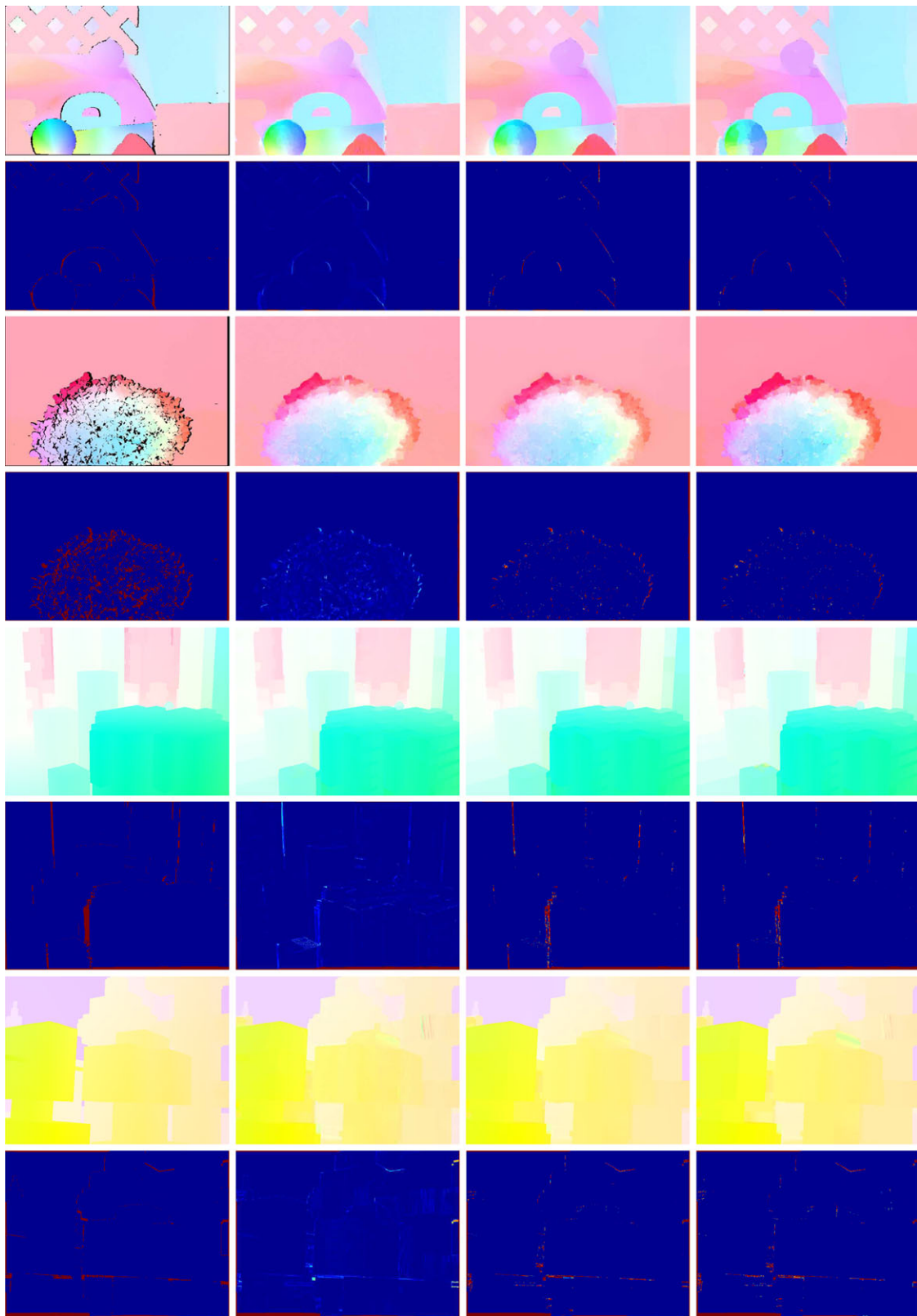


Fig. 6 This figure presents the occlusion and motion estimates on the sequences RubberWhale, Hydrangea, Urban2 and Urban3 from Middlebury dataset. Each sequence occupies two rows. The *odd rows*, left to right: ground truth optical flow, flow estimates before re-weighting stage and flow estimates after reweighting with Nesterov's algorithm

and split-Bregman method. The *even rows*, left to right: ground truth occluded regions, the initial estimate of the error term e_1 , the estimate of $|We_1|$ after the reweighting step with Nesterov's algorithm and split-Bregman method. The parts that are carried outside of the image domain under \hat{v} are also marked on the occlusion maps

Table 1 A comparison of the F-measure of our algorithm and Kolmogorov and Zabih (2001) on the Middlebury dataset. Since Kolmogorov and Zabih (2001) provide an occlusion detector whose output is binary, we simply compute the precision and recall of the output and report the F-measure based on these values. For comparison, we chose non zero elements of e_1 as detected occlusions and provide F-measure with respect to them

	Venus	RubberWhale	Hydrangea	Grove2	Grove3	Urban2	Urban3
F-measure (Kolmogorov and Zabih 2001)	0.63	0.28	0.31	0.62	0.52	0.43	0.53
F-measure (our method)	0.77	0.52	0.37	0.67	0.60	0.69	0.83

Table 2 The comparison of convergence time of the split-Bregman method and Nesterov's algorithm

	Venus	RubberWhale	Hydrangea	Grove2	Grove3	Urban2	Urban3
Nesterov's algorithm	222 s	342 s	355 s	463 s	494 s	499 s	483 s
Split Bregman method	90 s	111 s	133 s	260 s	360 s	288 s	277 s

Table 3 Quantitative comparison of the proposed models and other robust flow estimation methods (Black and Anandan 1996; Sun et al. 2010; Wedel et al. 2008) in terms of Average Angular Error (AAE)/Average End Point Error (AEPE)

	Venus	RubberWhale	Hydrangea	Grove2	Grove3	Urban2	Urban3
Huber- ℓ_1 -TV (Nesterov)	3.99/0.28	2.94/0.09	2.09/0.17	2.19/0.15	6.78/0.67	2.59/0.29	4.35/0.66
ℓ_2 -reweighted- ℓ_1 -TV (Nesterov)	3.96/0.31	5.09/0.16	2.36/0.19	2.60/0.17	7.71/0.78	3.41/0.38	4.91/0.76
ℓ_2 -reweighted- ℓ_1 -TV (split-Bregman)	4.09/0.33	4.91/0.16	2.34/0.19	2.31/0.15	7.72/0.75	2.86/0.35	4.20/0.63
Black & Anandan (Black and Anandan 1996)	7.81/0.44	5.06/0.14	2.48/0.21	2.76/0.20	6.90/0.75	4.06/0.54	11.18/0.94
Classic-L (Sun et al. 2010)	4.75/0.29	3.15/0.09	2.06/0.17	2.49/0.17	6.49/0.66	2.96/0.37	4.72/0.60
Wedel et al. (Improved L1-TV) (Wedel et al. 2008)	4.45/0.30	3.61/0.11	2.25/0.18	3.26/0.23	7.07/0.69	2.74/0.36	6.26/0.64

convergence speed, Figs. 4, 5 and 6, Tables 2 and 3. Both methods provide similar performance both in occlusion detection and motion estimation. However, split-Bregman method converges significantly faster, Table 2.

We have evaluated the accuracy of the flow estimates of our method and compared to other robust flow estimation techniques (Black and Anandan 1996; Sun et al. 2010; Wedel et al. 2008), Table 3. Huber- ℓ_1 -TV model minimized with Nesterov's algorithm provides superior accuracy. However, once the re-weighting stage is initialized with these estimates, and flow estimation is performed on the original images instead of the pre-filtered ones, the accuracy decreases.

Finally, we have evaluated the performance of our algorithm on the Middlebury evaluation dataset, Table 4. Nesterov's algorithm is used to estimate the optical flow on the evaluation set. In this experiment, we constructed the image pyramid with 11 levels with scale factor 0.75 and applied 10 warpings at each level. The flow fields estimated before the reweighting stage are ranked 10th and 12th in terms of AEPE and AAE respectively while the estimates after reweighting step are ranked 13th in terms of both error measures as of June 7, 2011. Our method also detects the occluded regions accurately on most of the sequences at evaluation set, Fig. 7.

6 Discussion

We have presented an algorithm to detect occlusions and establish correspondence between two images. It leverages on a formulation that, starting from standard assumptions (Lambertian reflection, constant illumination), arrives at a variational optimization problem. We have shown how this problem can be relaxed into a sequence of convex optimization schemes, each having a globally optimal solution, and presented two efficient numerical schemes for solving it.

We emphasize that our approach does *not* assume a rigid scene, or a single moving object. It also does *not* assume that the occluded region is simply connected. Instead, our model is general under the assumptions (a)–(b) as we show in Appendix A, and allows arbitrary (piece-wise diffeomorphic) domain deformations, corresponding to an arbitrary number of moving or deforming objects, and an arbitrary number of simply connected occluded regions (jointly represented by a multiply-connected domain Ω).

The fact that occlusion detection reduces to a two-phase segmentation of the domain into occluded Ω and visible region $D \setminus \Omega$ should not confuse the reader familiar with the *image segmentation* literature whereby two-phase segmentation of *one object* (foreground) from the background

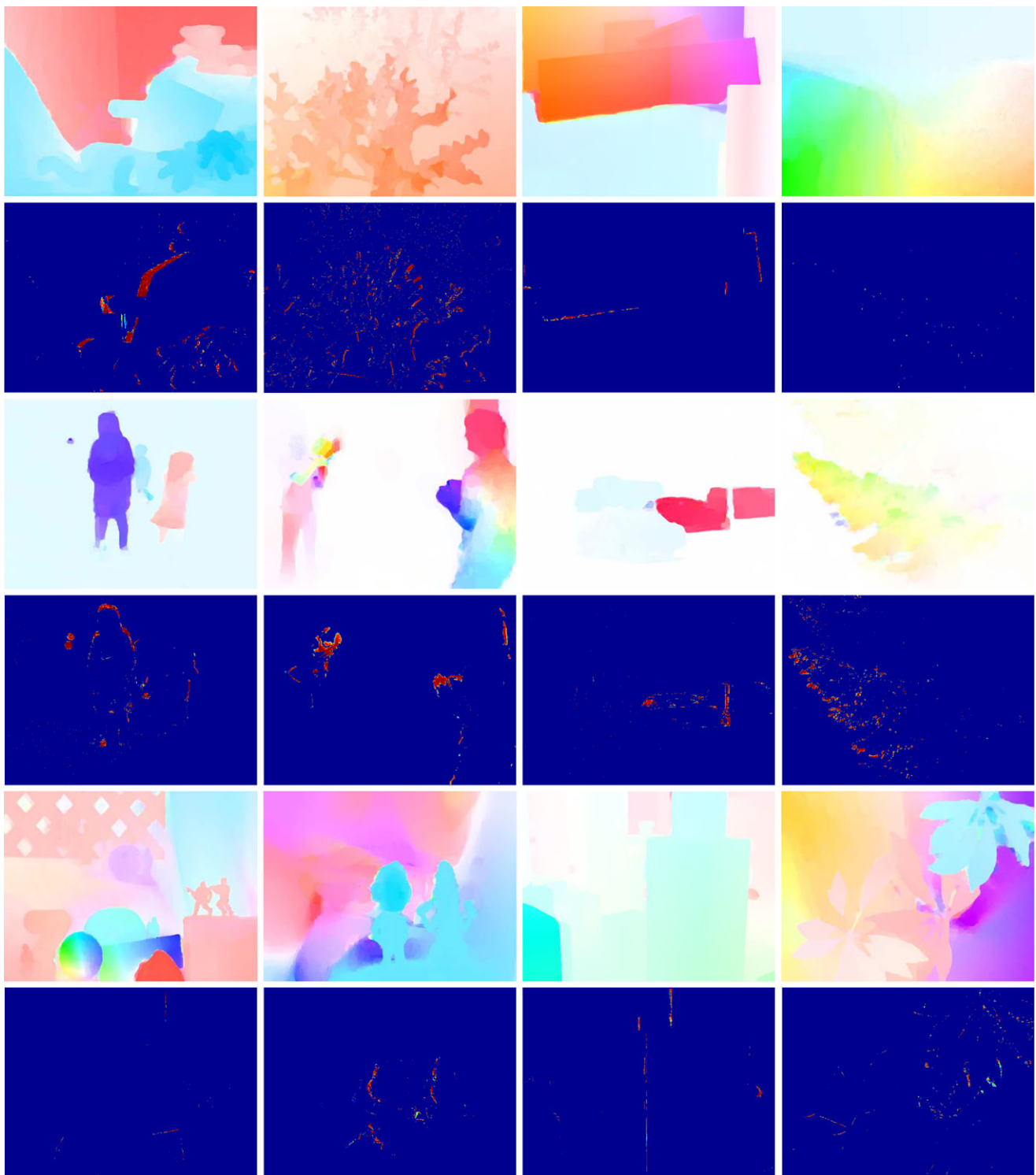


Fig. 7 Occlusion and flow estimates on the sequences Teddy, Grove, Wooden, Yosemite, Backyard, Basketball, Dumptruck, Evergreen, Army, Mequon, Urban and Schefflera from Middlebury evaluation set:

Estimated flow fields are depicted at *odd rows* while the estimates of $|We_1|$ after the reweighting step are illustrated at *even ones*

can be posed as a convex optimization problem (Chan et al. 2006). Note that in the approach of Chan et al. (2006) the problem can be made convex only in the occluded region

term, e_1 , but not jointly in e_1 and the motion field, v . Therefore, such an approach does not in general yield a global minimum.

Table 4 Quantitative evaluation of the proposed models on Middlebury test set in terms of AAE/AEPE

	Army	Mequon	Schefflera	Wooden	Grove	Urban	Yosemite	Teddy
w/o reweighting	3.62/0.09	2.92/0.16	4.49/0.22	3.14/0.17	3.26/0.75	3.52/0.34	5.10/0.22	2.02/0.53
w reweighting	4.30/0.11	4.03/0.30	5.13/0.42	3.56/0.20	3.23/0.75	3.22/0.33	2.81/0.15	1.94/0.55

The limitations of our approach stand mostly in its dependency from the regularization coefficients λ , μ and coefficient σ in the optimization. In the absence of some estimate of the variance coefficient λ , one is left with painstakingly tuning it by trial-and-error. Similarly, μ is a parameter that, like in any classification problem, trades off missed detections and false alarms, and therefore no single value is “optimal” in any meaningful sense. These limitations are shared by most variational optical flow estimation algorithms.

Acknowledgements Research supported by ARO 56765-CI, ONR N000141110863, AFOSR FA9550-09-1-0427.

Appendix A: Ambient-Lambert Model

In this section we show how to go from the assumptions (a)–(c) in Sect. 1 to (6). Let the scene $\{S, \rho\}$ be described by shape $S \subset \mathbb{R}^3$ (a collection of piece-wise smooth surfaces) and reflectance $\rho : S \rightarrow \mathbb{R}^k$ (diffuse albedo). Deviations from diffuse reflectance will not be modeled explicitly and lumped as error (inter-reflection, sub-surface scattering, specular reflection, cast shadows). Coarse illumination changes are modeled as a contrast transformation of the image range, and all other illumination effects are lumped into the additive error. The large number of independent phenomena being aggregated into such an error make it suitable to be modeled as a Gaussian random process ((2)-ii). Under these assumptions, the (isotropic) radiance ρ emitted by an area element around a point $p \in S$ is modulated by a monotonic continuous transformation m to yield the irradiance I measured at a pixel element x , except for the discrepancy $n : D \rightarrow \mathbb{R}_+^k$, and the correspondence between the point $p \in S$ and the pixel $x \in D$ is due to the motion of the viewer $g \in SE(3)$, the special Euclidean group of rotations and translations in three dimensional (3-D) space:

$$\begin{cases} I(x, t) = m(t) \circ \rho(p) \\ \quad + n(x, t); & p \in S \\ x = \pi(g(t)p); & x \in \pi(g(t)S) \\ I(x, t) = v(x, t) & x|g^{-1}(t)\pi^{-1}(x) \notin S \end{cases} \quad (30)$$

where $\pi : \mathbb{R}^3 \rightarrow \mathbb{R}^2$; $x \mapsto [x_1/x_3, x_2/x_3]^T$ is a central perspective projection. Away from the co-visible portion of the scene S , the image can take any value $v(x, t)$. Without loss of generality, the co-visible portion of the scene S

can be parametrized as the graph of a function (depth map), $\rho(x_0) = \bar{x}_0 Z(x_0)$, then the composition of maps

$$w : D \rightarrow \mathbb{R}^2; \quad x_0 \mapsto x = w(x_0) \doteq \pi(g\bar{x}_0 Z(x_0)) \quad (31)$$

spans the entire group of diffeomorphisms. This is the *motion field*, which is approximated by the optical flow when assumptions (a)–(c) are satisfied. Here a bar $\bar{x} \in \mathbb{P}^2$ denotes the homogeneous (projective) coordinates of the point with Euclidean coordinates $x \in \mathbb{R}^2$. Combining the two equations above, we have the two equivalent representations: $I(w(x_0)) = m \circ \rho(x_0) + n(w(x_0))$, $x_0 \in w^{-1}(D \setminus \Omega)$, or

$$I(x) = m \circ \rho(w^{-1}(x)) + n(x) \quad x \in D \setminus \Omega \quad (32)$$

with a slight abuse of notation since we have parametrized $\rho : S \rightarrow \mathbb{R}^k$ with one of the image planes, via $\rho(x) \leftarrow \rho(p(x))$, and we have re-defined $n(x) \leftarrow n(w^{-1}(x))$. Here D is the domain of the image, and Ω is the subset of the image where the object of interest is not visible (partial occlusion).

It can be shown that m can be eliminated via pre-processing by designing a representation that is a complete invariant statistic, that is a function of the image that is equivalent to it but for the effects of a contrast transformation (Caselles et al. 1999). There are several such functionals, including the curvature of the level sets of the image, or its dual (the gradient direction), or a normalization of contrast and offset of the image intensity, or spectral ratios if color images are available. In any case, we indicate this pre-processing via

$$\phi(I) = \phi(m \circ I) \quad (33)$$

Correspondence between two image regions can be established when they back-project onto the same portion of the scene S , or when that portion of the scene is *co-visible*. Therefore, establishing correspondence means, essentially, finding a scene (a shape S and an albedo ρ) that, under proper viewing conditions including a motion w and a contrast transformation h , yields a portion of each of the (two or more) images. This can be posed as an optimization problem, which under the assumptions (a)–(b) can be successively reduced into fewer and fewer unknowns:

$$\arg \min_{m, \rho, g, S} \int_{D \setminus \Omega} |I(x, t) - m \circ \rho \circ \pi(gS)| dx$$

(Theorem. 7.4, Robert 2001)

$$= \arg \min_{\rho, w} \int_{D \setminus \Omega} |\phi(I(x, t)) - \phi(\rho \circ w)| dx$$

(Theorem. 1, Soatto and Yezzi 2002)

$$= \arg \min_w \int_{D \setminus \Omega} |\phi(I(x, t)) - \phi(I(x, t - dt) \circ w)| dx$$

Of course, the (possibly multiply-connected) region Ω is also unknown, and can be represented via its characteristic function:

$$e_1(x) = \chi(\Omega) \quad (34)$$

where $\chi : D \rightarrow \mathbb{R}^+$ is such that $\chi(x) = 1$ if $x \in \Omega$, and $\chi(x) = 0$ elsewhere.

To ease the notational burden, we will assume that contrast has been eliminated via pre-processing, and drop the use of the function ϕ , so we re-define $I \leftarrow \phi(I)$. Writing explicitly the dependency of the “next image” on the occlusion domain, we have

$$\arg \min_w \int_{D \setminus \Omega} |I(x, t + dt) - I(w(x, t), t)|^2 dx \quad (35)$$

which is the \mathbb{L}^2 component of ψ_{data} in (6).

References

- Alvarez, L., Deriche, R., Papadopoulos, T., & Sánchez, J. (2007). Symmetrical dense optical flow estimation with occlusions detection. *International Journal of Computer Vision*, 75(3), 371–385.
- Ayvaci, A., & Soatto, S. (2011). Efficient model selection for detachable object detection. In *Proc. of energy minimization methods in computer vision and pattern recognition*, July 2011.
- Ayvaci, A., Raptis, M., & Soatto, S. (2010). Occlusion detection and motion estimation with convex optimization. In *Advances in neural information processing systems*.
- Baker, S., Scharstein, D., Lewis, J., Roth, S., Black, M., & Szeliski, R. (2007). A database and evaluation methodology for optical flow. In *Proc. of the international conference on computer vision* (pp. 1–8).
- Becker, S., Bobin, J., & Candes, E. (2009). NESTA: a fast and accurate first-order method for sparse recovery. Arxiv preprint arXiv, 904.
- Ben-Ari, R., & Sochen, N. (2007). Variational stereo vision with sharp discontinuities and occlusion handling. In *Proc. of international conference on computer vision* (pp. 1–7).
- Black, M., & Anandan, P. (1996). The robust estimation of multiple motions: Parametric and piecewise-smooth flow fields. *Computer Vision and Image Understanding*, 63(1), 75–104.
- Brox, T., Bruhn, A., Papenberg, N., & Weickert, J. (2004). High accuracy optical flow estimation based on a theory for warping. In *Proc. of European conference on computer vision* (pp. 25–36).
- Bruhn, A., Weickert, J., & Schnörr, C. (2005). Lucas/Kanade meets Horn/Schunck: Combining local and global optic flow methods. *International Journal of Computer Vision*, 61(3), 211–231.
- Candes, E., Wakin, M., & Boyd, S. (2008). Enhancing sparsity by reweighted L1 minimization. *The Journal of Fourier Analysis and Applications*, 14(5), 877–905.
- Caselles, V., Coll, B., & Morel, J.-M. (1999). Topographic maps and local contrast changes in natural images. *International Journal of Computer Vision*, 33(1), 5–27.
- Chan, T., Esedoglu, S., & Nikolova, M. (2006). Algorithms for finding global minimizers of denoising and segmentation models. *SIAM Journal on Applied Mathematics*, 66(1), 1632–1648.
- Dahl, J., Hansen, P., Jensen, S., & Jensen, T. (2009). Algorithms and software for total variation image reconstruction via first-order methods. *Numerical Algorithms*, 67–92.
- Gibson, J. J. (1984). *The ecological approach to visual perception*. LEA.
- Goldstein, T., & Osher, S. (2009). The split Bregman method for L1 regularized problems. *SIAM Journal on Imaging Sciences*, 2(2), 323–343.
- He, X., & Yuille, A. (2010). Occlusion boundary detection using pseudo-depth. In *Proc. of the European conference on computer vision*.
- Horn, B., & Schunck, B. (1981). Determining optical flow. *Computer Vision*, 17, 185–203.
- Humayun, A., Mac Aodha, O., & Brostow, G. J. (2011). Learning to find occlusion regions. In *Proc. of conference on computer vision and pattern recognition*.
- Ince, S., & Konrad, J. (2008). Occlusion-aware optical flow estimation. *IEEE Transactions on Image Processing*, 17(8), 1443–1451.
- Jackson, J. D., Yezzi, A. J., & Soatto, S. (2005). Dynamic shape and appearance modeling via moving and deforming layers. In *Proc. of workshop on energy minimization in computer vision and pattern recognition (EMMCVPR)* (pp. 427–438).
- Jackson, J., Yezzi, A. J., & Soatto, S. (2008). Dynamic shape and appearance modeling via moving and deforming layers. *International Journal of Computer Vision*.
- Kim, Y., Martínez, A., & Kak, A. (2005). Robust motion estimation under varying illumination. *Image and Vision Computing*, 23(4), 365–375.
- Kolmogorov, V., & Zabih, R. (2001). Computing visual correspondence with occlusions via graph cuts. In *Proc. of international conference on computer vision* (pp. 508–515).
- Lim, K., Das, A., & Chong, M. (2002). Estimation of occlusion and dense motion fields in a bidirectional Bayesian framework. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 712–718.
- Negahdaripour, S. (1998). Revised definition of optical flow: integration of radiometric and geometric cues for dynamic scene analysis. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 961–979.
- Nesterov, Y. (1983). A method for unconstrained convex minimization problem with the rate of convergence $O(1/k^2)$. *Doklady Akademii Nauk SSSR*, 269, 543–547.
- Nesterov, Y. (2005). Smooth minimization of non-smooth functions. *Mathematical Programming*, 103(1), 127–152.
- Proesmans, M., Van Gool, L., & Oosterlinck, A. (1994). Determination of optical flow and its discontinuities using a non-linear diffusion. In *Proc. of European conference of computer vision*.
- Robert, C. P. (2001). *The Bayesian choice*. New York: Springer.
- Rudin, L., Osher, S., & Fatemi, E. (1992). Nonlinear total variation based noise removal algorithms. *Physica. D*, 60, 259–268.
- Shulman, D., & Herve, J. (1989). Regularization of discontinuous flow fields. In *Proc. of workshop on visual motion* (pp. 81–86).
- Soatto, S. (2011). Actionable information in vision. In *Machine learning for computer vision*. Berlin: Springer.
- Soatto, S., & Yezzi, A. (2002). Deformation: deforming motion, shape average and the joint segmentation and registration of images. In *Proc. of the European conference on computer vision* (Vol. 3, pp. 32–47).
- Soatto, S., Yezzi, A. J., & Jin, H. (2003). Tales of shape and radiance in multiview stereo. In *Proc. of international conference on computer vision* (pp. 974–981). October 2003.

- Stein, A., & Hebert, M. (2009). Occlusion boundaries from motion: low-level detection and mid-level reasoning. *International Journal of Computer Vision*, 82(3), 325–357.
- Strecha, C., Fransens, R., & Van Gool, L. (2004). A probabilistic approach to large displacement optical flow and occlusion detection. In *ECCV workshop SMVP* (pp. 71–82). Berlin: Springer.
- Sun, J., Li, Y., Kang, S., & Shum, H. (2005). Symmetric stereo matching for occlusion handling. In *Proc. of conference on computer vision and pattern recognition* (Vol. 2, p. 399).
- Sun, D., Roth, S., & Black, M. (2010). Secrets of optical flow estimation and their principles. In *Proc. of conference on computer vision and pattern recognition* (pp. 2432–2439).
- Sundaramoorthi, G., Petersen, P., Varadarajan, V. S., & Soatto, S. (2009). On the set of images modulo viewpoint and contrast changes. In *Proc. of conference on computer vision and pattern recognition*.
- Teng, C., Lai, S., Chen, Y., & Hsu, W. (2005). Accurate optical flow computation under non-uniform brightness variations. *Computer Vision and Image Understanding*, 97(3), 315–346.
- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B. Methodological*, 58(1), 267–288.
- Verri, A., & Poggio, T. (1989). Motion field and optical flow: Qualitative properties. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 11(5), 490–498.
- Wang, J., & Adelson, E. (1994). Representing moving images with layers. *IEEE Transactions on Image Processing*, 3(5), 625–638.
- Wang, Y., Yin, W., & Zhang, Y. (2007). *A fast algorithm for image deblurring with total variation regularization* (CAAM Technical Reports). Rice University.
- Wedel, A., Pock, T., Zach, C., Bischof, H., & Cremers, D. (2008). An improved algorithm for TV-L1 optical flow. In *Proc. of statistical and geometrical approaches to visual motion analysis: International Dagstuhl seminar*.
- Wedel, A., Cremers, D., Pock, T., & Bischof, H. (2009). Structure- and motion-adaptive regularization for high accuracy optic flow. In *Proc. of international conference on computer vision*.
- Weiss, P., Blanc-Féraud, L., & Aubert, G. (2009). Efficient schemes for total variation minimization under constraints in image processing. *SIAM Journal on Scientific Computing*, 31, 2047.
- Werlberger, M., Trobin, W., Pock, T., Wedel, A., Cremers, D., & Bischof, H. (2009). Anisotropic Huber-L1 optical flow. In *Proc. of British machine vision conference*.
- Xiao, J., Cheng, H., Sawhney, H., Rao, C., Isnardi, M., et al. (2006). Bilateral filtering-based optical flow estimation with occlusion detection. In *Proc. of European conference of computer vision* (Vol. 3951, p. 211).