

A Database and Evaluation Methodology for Optical Flow

Simon Baker · Daniel Scharstein · J.P. Lewis ·
Stefan Roth · Michael J. Black · Richard Szeliski

Received: 18 December 2009 / Accepted: 20 September 2010

© Springer Science+Business Media, LLC 2010. This article is published with open access at Springerlink.com

Abstract The quantitative evaluation of optical flow algorithms by Barron et al. (1994) led to significant advances in performance. The challenges for optical flow algorithms today go beyond the datasets and evaluation methods proposed in that paper. Instead, they center on problems associated with complex natural scenes, including nonrigid motion, real sensor noise, and motion discontinuities. We propose a new set of benchmarks and evaluation methods for the next generation of optical flow algorithms. To that end, we contribute four types of data to test different aspects of optical flow algorithms: (1) sequences with nonrigid motion where the ground-truth flow is determined by

tracking hidden fluorescent texture, (2) realistic synthetic sequences, (3) high frame-rate video used to study interpolation error, and (4) modified stereo sequences of static scenes. In addition to the average angular error used by Barron et al., we compute the absolute flow endpoint error, measures for frame interpolation error, improved statistics, and results at motion discontinuities and in textureless regions. In October 2007, we published the performance of several well-known methods on a preliminary version of our data to establish the current state of the art. We also made the data freely available on the web at <http://vision.middlebury.edu/flow/>. Subsequently a number of researchers have uploaded their results to our website and published papers using the data. A significant improvement in performance has already been achieved. In this paper we analyze the results obtained to date and draw a large number of conclusions from them.

A preliminary version of this paper appeared in the IEEE International Conference on Computer Vision (Baker et al. 2007).

S. Baker · R. Szeliski
Microsoft Research, Redmond, WA, USA

S. Baker
e-mail: sbaker@microsoft.com

R. Szeliski
e-mail: szeliski@microsoft.com

D. Scharstein (✉)
Middlebury College, Middlebury, VT, USA
e-mail: schar@middlebury.edu

J.P. Lewis
Weta Digital, Wellington, New Zealand
e-mail: zilla@computer.org

S. Roth
TU Darmstadt, Darmstadt, Germany
e-mail: sroth@cs.tu-darmstadt.de

M.J. Black
Brown University, Providence, RI, USA
e-mail: black@cs.brown.edu

Keywords Optical flow · Survey · Algorithms · Database · Benchmarks · Evaluation · Metrics

1 Introduction

As a subfield of computer vision matures, datasets for quantitatively evaluating algorithms are essential to ensure continued progress. Many areas of computer vision, such as stereo (Scharstein and Szeliski 2002), face recognition (Philips et al. 2005; Sim et al. 2003; Gross et al. 2008; Georgiades et al. 2001), and object recognition (Fei-Fei et al. 2006; Everingham et al. 2009), have challenging datasets to track the progress made by leading algorithms and to stimulate new ideas. Optical flow was actually one of the first areas to have such a benchmark, introduced by Barron et al. (1994). The field benefited greatly from this

study, which led to rapid and measurable progress. To continue the rapid progress, new and more challenging datasets are needed to push the limits of current technology, reveal where current algorithms fail, and evaluate the next generation of optical flow algorithms. Such an evaluation dataset for optical flow should ideally consist of complex real scenes with all the artifacts of real sensors (noise, motion blur, etc.). It should also contain substantial motion discontinuities and nonrigid motion. Of course, the image data should be paired with dense, subpixel-accurate, ground-truth flow fields.

The presence of nonrigid or independent motion makes collecting a ground-truth dataset for optical flow far harder than for stereo, say, where structured light (Scharstein and Szeliski 2002) or range scanning (Seitz et al. 2006) can be used to obtain ground truth. Our solution is to collect four different datasets, each satisfying a different subset of the desirable properties above. The combination of these datasets provides a basis for a thorough evaluation of current optical flow algorithms. Moreover, the relative performance of algorithms on the different datatypes may stimulate further research. In particular, we collected the following four types of data:

- *Real Imagery of Nonrigidly Moving Scenes:* Dense ground-truth flow is obtained using hidden fluorescent texture painted on the scene. We slowly move the scene, at each point capturing separate test images (in visible light) and ground-truth images with trackable texture (in UV light). Note that a related technique is being used commercially for motion capture (Mova LLC 2004) and Tappen et al. (2006) recently used certain wavelengths to hide ground truth in intrinsic images. Another form of hidden markers was also used in Ramnath et al. (2008) to provide a sparse ground-truth alignment (or flow) of face images. Finally, Liu et al. recently proposed a method to obtain ground-truth using human annotation (Liu et al. 2008).
- *Realistic Synthetic Imagery:* We address the limitations of simple synthetic sequences such as *Yosemite* (Barron et al. 1994) by rendering more complex scenes with larger motion ranges, more realistic texture, independent motion, and with more complex occlusions.
- *Imagery for Frame Interpolation:* Intermediate frames are withheld and used as ground truth. In a wide class of applications such as video re-timing, novel-view generation, and motion-compensated compression, what is important is not how well the flow matches the ground-truth motion, but how well intermediate frames can be predicted using the flow (Szeliski 1999).
- *Real Stereo Imagery of Rigid Scenes:* Dense ground truth is captured using structured light (Scharstein and Szeliski 2003). The data is then adapted to be more appropriate for optical flow by cropping to make the disparity range roughly symmetric.

We collected enough data to be able to split our collection into a training set (12 datasets) and a final evaluation set (12 datasets). The training set includes the ground truth and is meant to be used for debugging, parameter estimation, and possibly even learning (Sun et al. 2008; Li and Huttenlocher 2008). The ground truth for the final evaluation set is not publicly available (with the exception of the *Yosemite* sequence, which is included in the test set to allow some comparison with algorithms published prior to the release of our data).

We also extend the set of performance measures and the evaluation methodology of Barron et al. (1994) to focus attention on current algorithmic problems:

- *Error Metrics:* We report both average angular error (Barron et al. 1994) and flow endpoint error (pixel distance) (Otte and Nagel 1994). For image interpolation, we compute the residual RMS error between the interpolated image and the ground-truth image. We also report a gradient-normalized RMS error (Szeliski 1999).
- *Statistics:* In addition to computing averages and standard deviations as in Barron et al. (1994), we also compute robustness measures (Scharstein and Szeliski 2002) and percentile-based accuracy measures (Seitz et al. 2006).
- *Region Masks:* Following Scharstein and Szeliski (2002), we compute the error measures and their statistics over certain masked regions of research interest. In particular, we compute the statistics near motion discontinuities and in textureless regions.

Note that we require flow algorithms to estimate a dense flow field. An alternate approach might be to allow algorithms to provide a confidence map, or even to return a sparse or incomplete flow field. Scoring such outputs is problematic, however. Instead, we expect algorithms to generate a flow estimate everywhere (for instance, using internal confidence measures to fill in areas with uncertain flow estimates due to lack of texture).

In October 2007 we published the performance of several well-known algorithms on a preliminary version of our data to establish the current state of the art (Baker et al. 2007). We also made the data freely available on the web at <http://vision.middlebury.edu/flow/>. Subsequently a large number of researchers have uploaded their results to our website and published papers using the data. A significant improvement in performance has already been achieved. In this paper we present both results obtained by classic algorithms, as well as results obtained since publication of our preliminary data. In addition to summarizing the overall conclusions of the currently uploaded results, we also examine how the results vary: (1) across the metrics, statistics, and region masks, (2) across the various datatypes and datasets, (3) from flow estimation to interpolation, and (4) depending on the components of the algorithms.

The remainder of this paper is organized as follows. We begin in Sect. 2 with a survey of existing optical flow algorithms, benchmark databases, and evaluations. In Sect. 3 we describe the design and collection of our database, and briefly discuss the pros and cons of each dataset. In Sect. 4 we describe the evaluation metrics. In Sect. 5 we present the experimental results and discuss the major conclusions that can be drawn from them.

2 Related Work and Taxonomy of Optical Flow Algorithms

Optical flow estimation is an extensive field. A fully comprehensive survey is beyond the scope of this paper. In this related work section, our goals are: (1) to present a taxonomy of the main components in the majority of existing optical flow algorithms, and (2) to focus primarily on recent work and place the contributions of this work in the context of our taxonomy. Note that our taxonomy is similar to those of Stiller and Konrad (1999) for optical flow and Scharstein and Szeliski (2002) for stereo. For more extensive coverage of older work, the reader is referred to previous surveys such as those by Aggarwal and Nandhakumar (1988), Barron et al. (1994), Otte and Nagel (1994), Mitiche and Bouthemy (1996), and Stiller and Konrad (1999).

We first define what we mean by optical flow. Following Horn's (1986) taxonomy, the *motion field* is the 2D projection of the 3D motion of surfaces in the world, whereas the *optical flow* is the *apparent motion* of the brightness patterns in the image. These two motions are not always the same and, in practice, the goal of 2D motion estimation is application dependent. In frame interpolation, it is preferable to estimate apparent motion so that, for example, specular highlights move in a realistic way. On the other hand, in applications where the motion is used to interpret or reconstruct the 3D world, the motion field is what is desired.

In this paper, we consider both motion field estimation and apparent motion estimation, referring to them collectively as optical flow. The ground truth for most of our datasets is the true motion field, and hence this is how we define and evaluate optical flow accuracy. For our interpolation datasets, the ground truth consists of images captured at an intermediate time instant. For this data, our definition of optical flow is really the apparent motion.

We do, however, restrict attention to optical flow algorithms that estimate a separate 2D motion vector for each pixel in one frame of a sequence or video containing two or more frames. We exclude transparency which requires multiple motions per pixel. We also exclude more global representations of the motion such as parametric motion estimates (Bergen et al. 1992).

Most existing optical flow algorithms pose the problem as the optimization of a global energy function that is the weighted sum of two terms:

$$E_{\text{Global}} = E_{\text{Data}} + \lambda E_{\text{Prior}}. \quad (1)$$

The first term E_{Data} is the *Data Term*, which measures how consistent the optical flow is with the input images. We consider the choice of the data term in Sect. 2.1. The second term E_{Prior} is the *Prior Term*, which favors certain flow fields over others (for example E_{Prior} often favors smoothly varying flow fields). We consider the choice of the prior term in Sect. 2.2. The optical flow is then computed by optimizing the global energy E_{Global} . We consider the choice of the optimization algorithm in Sects. 2.3 and 2.4. In Sect. 2.5 we consider a number of miscellaneous issues. Finally, in Sect. 2.6 we survey previous databases and evaluations.

2.1 Data Term

2.1.1 Brightness Constancy

The basis of the data term used by most algorithms is *Brightness Constancy*, the assumption that when a pixel flows from one image to another, its intensity or color does not change. This assumption combines a number of assumptions about the reflectance properties of the scene (e.g., that it is Lambertian), the illumination in the scene (e.g., that it is uniform—Vedula et al. 2005) and about the image formation process in the camera (e.g., that there is no vignetting). If $I(x, y, t)$ is the intensity of a pixel (x, y) at time t and the flow is $(u(x, y, t), v(x, y, t))$, Brightness Constancy can be written as:

$$I(x, y, t) = I(x + u, y + v, t + 1). \quad (2)$$

Linearizing (2) by applying a first-order Taylor expansion to the right-hand side yields the approximation:

$$I(x, y, t) = I(x, y, t) + u \frac{\partial I}{\partial x} + v \frac{\partial I}{\partial y} + 1 \frac{\partial I}{\partial t}, \quad (3)$$

which simplifies to the *Optical Flow Constraint* equation:

$$u \frac{\partial I}{\partial x} + v \frac{\partial I}{\partial y} + \frac{\partial I}{\partial t} = 0. \quad (4)$$

Both Brightness Constancy and the Optical Flow Constraint equation provide just one constraint on the two unknowns at each pixel. This is the origin of the *Aperture Problem* and the reason that optical flow is ill-posed and must be regularized with a prior term (see Sect. 2.2).

The data term E_{Data} can be based on either Brightness Constancy in (2) or on the Optical Flow Constraint in (4). In either case, the equation is turned into an error per pixel,

the set of which is then aggregated over the image in some manner (see Sect. 2.1.2). If Brightness Constancy is used, it is generally converted to the Optical Flow Constraint during the derivation of most continuous optimization algorithms (see Sect. 2.3), which often involves the use of a Taylor expansion to linearize the energies. The two constraints are therefore essentially equivalent in practical algorithms (Brox et al. 2004).

An alternative to the assumption of “constancy” is that the signals (images) at times t and $t + 1$ are highly *correlated* (Pratt 1974; Burt et al. 1982). Various correlation constraints can be used for computing dense flow including normalized cross correlation and Laplacian correlation (Burt et al. 1983; Glazer et al. 1983; Sun 1999).

2.1.2 Choice of the Penalty Function

Equations (2) and (4) both provide one error per pixel, which leads to the question of how these errors are aggregated over the image. A baseline approach is to use an L2 norm as in the Horn and Schunck algorithm (Horn and Schunck 1981):

$$E_{\text{Data}} = \sum_{x,y} \left[u \frac{\partial I}{\partial x} + v \frac{\partial I}{\partial y} + \frac{\partial I}{\partial t} \right]^2. \quad (5)$$

If (5) is interpreted probabilistically, the use of the L2 norm means that the errors in the Optical Flow Constraint are assumed to be Gaussian and IID. This assumption is rarely true in practice, particularly near occlusion boundaries where pixels at time t may not be visible at time $t + 1$. Black and Anandan (1996) present an algorithm that can use an arbitrary robust penalty function, illustrating their approach with the specific choice of a Lorentzian penalty function. A common choice by a number of recent algorithms (Brox et al. 2004; Wedel et al. 2008) is the L1 norm, which is sometimes approximated with a differentiable version:

$$\|\mathbf{E}\|_1 = \sum_{x,y} |E_{x,y}| \approx \sum_{x,y} \sqrt{\|E_{x,y}\|^2 + \epsilon^2}, \quad (6)$$

where \mathbf{E} is a vector of errors $E_{x,y}$, $\|\cdot\|_1$ denotes the L1 norm, and ϵ is a small positive constant. A variety of other penalty functions have been used.

2.1.3 Photometrically Invariant Features

Instead of using the raw intensity or color values in the images, it is also possible to use features computed from those images. In fact, some of the earliest optical flow algorithms used filtered images to reduce the effects of shadows (Burt et al. 1983; Anandan 1989). One recently popular choice (for example used in Brox et al. 2004 among others) is to augment or replace (2) with a similar term based on the gradient of the image:

$$\nabla I(x, y, t) = \nabla I(x + u, y + v, t + 1). \quad (7)$$

Empirically the gradient is often more robust to (approximately additive) illumination changes than the raw intensities. Note, however, that (7) makes the additional assumption that the flow is locally translational; e.g., local scale changes, rotations, etc., can violate (7) even when (2) holds. It is also possible to use more complicated features than the gradient. For example a Field-of-Experts formulation is used in Sun et al. (2008) and SIFT features are used in Liu et al. (2008).

2.1.4 Modeling Illumination, Blur, and Other Appearance Changes

The motivation for using features is to increase robustness to illumination and other appearance changes. Another approach is to estimate the change explicitly. For example, suppose $g(x, y)$ denotes a multiplicative scale factor and $b(x, y)$ an additive term that together model the illumination change between $I(x, y, t)$ and $I(x, y, t + 1)$. Brightness Constancy in (2) can be generalized to:

$$g(x, y)I(x, y, t) = I(x + u, y + v, t + 1) + b(x, y). \quad (8)$$

Note that putting $g(x, y)$ on the left-hand side is preferable to putting it on the right-hand side as it can make optimization easier (Seitz and Baker 2009). Equation (8) is even more under-constrained than (2), with four unknowns per pixel rather than two. It can, however, be solved by putting an appropriate prior on the two components of the illumination change model $g(x, y)$ and $b(x, y)$ (Negahdaripour 1998; Seitz and Baker 2009). Explicit illumination modeling can be generalized in several ways, for example to model the changes physically over a longer time interval (Haussecker and Fleet 2000) or to model blur (Seitz and Baker 2009).

2.1.5 Color and Multi-Band Images

Another issue, addressed by a number of authors (Ohta 1989; Markandey and Flinchbaugh 1990; Golland and Bruckstein 1997), is how to modify the data term for color or multi-band images. The simplest approach is to add a data term for each band, for example performing the summation in (5) over the color bands, as well as the pixel coordinates x, y . More sophisticated approaches include using the HSV color space and treating the bands differently (e.g., by using different weights or norms) (Zimmer et al. 2009).

2.2 Prior Term

The data term alone is ill-posed with fewer constraints than unknowns. It is therefore necessary to add a prior to favor one possible solution over another. Generally speaking, while most priors are smoothness priors, a wide variety of choices are possible.

2.2.1 First Order

Arguably the simplest prior is to favor small first-order derivatives (gradients) of the flow field. If we use an L2 norm, then we might, for example, define:

$$E_{\text{Prior}} = \sum_{x,y} \left[\left(\frac{\partial u}{\partial x} \right)^2 + \left(\frac{\partial u}{\partial y} \right)^2 + \left(\frac{\partial v}{\partial x} \right)^2 + \left(\frac{\partial v}{\partial y} \right)^2 \right]. \quad (9)$$

The combination of (5) and (9) defines the energy used by Horn and Schunck (1981). Given more than two frames in the video, it is also possible to add temporal smoothness terms $\frac{\partial u}{\partial t}$ and $\frac{\partial v}{\partial t}$ to (9) (Murray and Buxton 1987; Black and Anandan 1991; Brox et al. 2004). Note, however, that the temporal terms need to be weighted differently from the spatial ones.

2.2.2 Choice of the Penalty Function

As for the data term in Sect. 2.1.2, under a probabilistic interpretation, the use of an L2 norm assumes that the gradients of the flow field are Gaussian and IID. Again, this assumption is violated in practice and so a wide variety of other penalty functions have been used. The algorithm by Black and Anandan (1996) also uses a first-order prior, but can use an arbitrary robust penalty function on the prior term rather than the L2 norm in (9). While Black and Anandan (1996) use the same Lorentzian penalty function for both the data and spatial term, there is no need for them to be the same. The L1 norm is also a popular choice of penalty function (Brox et al. 2004; Wedel et al. 2008). When the L1 norm is used to penalize the gradients of the flow field, the formulation falls in the class of Total Variation (TV) methods.

There are two common ways such robust penalty functions are used. One approach is to apply the penalty function separately to each derivative and then to sum up the results. The other approach is to first sum up the squares (or absolute values) of the gradients and then apply a single robust penalty function. Some algorithms use the first approach (Black and Anandan 1996), while others use the second (Bruhn et al. 2005; Brox et al. 2004; Wedel et al. 2008).

Note that some penalty (log probability) functions have probabilistic interpretations related to the distribution of flow derivatives (Roth and Black 2007).

2.2.3 Spatial Weighting

One popular refinement for the prior term is one that weights the penalty function with a spatially varying function. One particular example is to vary the weight depending on the

gradient of the image:

$$E_{\text{Prior}} = \sum_{x,y} w(\nabla I) \left[\left(\frac{\partial u}{\partial x} \right)^2 + \left(\frac{\partial u}{\partial y} \right)^2 + \left(\frac{\partial v}{\partial x} \right)^2 + \left(\frac{\partial v}{\partial y} \right)^2 \right]. \quad (10)$$

Equation (10) could be used to reduce the weight of the prior at edges (high $|\nabla I|$) because there is a greater likelihood of a flow discontinuity at an intensity edge than inside a smooth region. The weight can also be a function of an over-segmentation of the image, rather than the gradient, for example down-weighting the prior between different segments (Seitz and Baker 2009).

2.2.4 Anisotropic Smoothness

In (10) the weighting function is isotropic, treating all directions equally. A variety of approaches weight the smoothness prior anisotropically. For example, Nagel and Enkelmann (1986) and Werlberger et al. (2009) weight the direction along the image gradient less than the direction orthogonal to it, and Sun et al. (2008) learn a Steerable Random Field to define the weighting. Zimmer et al. (2009) perform a similar anisotropic weighting, but the directions are defined by the data constraint rather than the image gradient.

2.2.5 Higher-Order Priors

The first-order priors in Sect. 2.2.1 can be replaced with priors that encourage the second-order derivatives ($\frac{\partial^2 u}{\partial x^2}$, $\frac{\partial^2 u}{\partial y^2}$, $\frac{\partial^2 u}{\partial x \partial y}$, $\frac{\partial^2 v}{\partial x^2}$, $\frac{\partial^2 v}{\partial y^2}$, $\frac{\partial^2 v}{\partial x \partial y}$) to be small (Anandan and Weiss 1985; Trobin et al. 2008).

A related approach is to use an affine prior (Ju et al. 1996; Ju 1998; Nir et al. 2008; Seitz and Baker 2009). One approach is to over-parameterize the flow (Nir et al. 2008). Instead of solving for two flow vectors $(u(x, y, t), v(x, y, t))$ at each pixel, the algorithm in Nir et al. (2008) solves for 6 affine parameters $a_i(x, y, t)$, $i = 1, \dots, 6$ where the flow is given by:

$$u(x, y, t) = a_1(x, y, t) + \frac{x - x_0}{x_0} a_3(x, y, t) + \frac{y - y_0}{y_0} a_5(x, y, t), \quad (11)$$

$$v(x, y, t) = a_2(x, y, t) + \frac{x - x_0}{x_0} a_4(x, y, t) + \frac{y - y_0}{y_0} a_6(x, y, t), \quad (12)$$

where (x_0, y_0) is the middle of the image. Equations (11) and (12) are then substituted into any of the data terms

above. Ju et al. formulate the prior so that neighboring affine parameters should be similar (Ju et al. 1996). As above, a robust penalty may be used and, further, may vary depending on the affine parameter (for example weighting a_1 and a_2 differently from $a_3 \dots a_6$).

2.2.6 Rigidity Priors

A number of authors have explored rigidity or fundamental matrix priors which, in the absence of other evidence, favor flows that are aligned with epipolar lines. These constraints have both been strictly enforced (Adiv 1985; Hanna 1991; Nir et al. 2008) and added as a soft prior (Wedel et al. 2008; Wedel et al. 2009; Valgaerts et al. 2008).

2.3 Continuous Optimization Algorithms

The two most commonly used continuous optimization techniques in optical flow are: (1) gradient descent algorithms (Sect. 2.3.1) and (2) extremal or variational approaches (Sect. 2.3.2). In Sect. 2.3.3 we describe a small number of other approaches.

2.3.1 Gradient Descent Algorithms

Let \mathbf{f} be a vector resulting from concatenating the horizontal and vertical components of the flow at every pixel. The goal is then to optimize E_{Global} with respect to \mathbf{f} . The simplest gradient descent algorithm is *steepest descent* (Baker and Matthews 2004), which takes steps in the direction of the negative gradient $-\frac{\partial E_{\text{Global}}}{\partial \mathbf{f}}$. An important question with steepest descent is how big the step size should be. One approach is to adjust the step size iteratively, increasing it if the algorithm makes a step that reduces the energy and decreasing it if the algorithm tries to make a step that increases the error. Another approach used in Black and Anandan (1996) is to set the step size to be:

$$-w \frac{1}{T} \frac{\partial E_{\text{Global}}}{\partial \mathbf{f}}. \quad (13)$$

In this expression, T is an upper bound on the second derivatives of the energy; $T \geq \frac{\partial^2 E_{\text{Global}}}{\partial f_i^2}$ for all components f_i in the vector \mathbf{f} . The parameter $0 < w < 2$ is an over-relaxation parameter. Without it, (13) tends to take too small steps because: (1) T is an upper bound, and (2) the equation does not model the off-diagonal elements in the Hessian. It can be shown that if E_{Global} is a quadratic energy function (i.e., the problem is equivalent to solving a large linear system), convergence to the global minimum can be guaranteed (albeit possibly slowly) for any $0 < w < 2$. In general E_{Global} is nonlinear and so there is no such guarantee. However, based on the theoretical result in the linear case, a value

around $w \approx 1.95$ is generally used. Also note that many non-quadratic (e.g., robust) formulations can be solved with iteratively reweighted least squares (IRLS); i.e., they are posed as a sequence of quadratic optimization problems with a data-dependent weighting function that varies from iteration to iteration. The weighted quadratic is iteratively solved and the weights re-estimated.

In general, steepest descent algorithms are relatively weak optimizers requiring a large number of iterations because they fail to model the coupling between the unknowns. A second-order model of this coupling is contained in the Hessian matrix $\frac{\partial^2 E_{\text{Global}}}{\partial f_i \partial f_j}$. Algorithms that use the Hessian matrix or approximations to it such as the Newton method, Quasi-Newton methods, the Gauss-Newton method, and the Levenberg-Marquardt algorithm (Baker and Matthews 2004) all converge far faster. These algorithms are however inapplicable to the general optical flow problem because they require estimating and inverting the Hessian, a $2n \times 2n$ matrix where there are n pixels in the image. These algorithms are applicable to problems with fewer parameters such as the Lucas-Kanade algorithm (Lucas and Kanade 1981) and variants (Le Besnerais and Champagnat 2005), which solve for a single flow vector (2 unknowns) independently for each block of pixels. Another set of examples are parametric motion algorithms (Bergen et al. 1992), which also just solve for a small number of unknowns.

2.3.2 Variational and Other Extremal Approaches

The second class of algorithms assume that the global energy function can be written in the form:

$$E_{\text{Global}} = \iint E(u(x, y), v(x, y), x, y, u_x, u_y, v_x, v_y) dx dy, \quad (14)$$

where $u_x = \frac{\partial u}{\partial x}$, $u_y = \frac{\partial u}{\partial y}$, $v_x = \frac{\partial v}{\partial x}$, and $v_y = \frac{\partial v}{\partial y}$. At this stage, $u = u(x, y)$ and $v = v(x, y)$ are treated as unknown 2D functions rather than the set of unknown parameters (the flows at each pixel). The parameterization of these functions occurs later. Note that (14) imposes limitations on the functional form of the energy, i.e., that it is just a function of the flow u, v , the spatial coordinates x, y and the gradients of the flow u_x, u_y, v_x and v_y . A wide variety of energy functions do satisfy this requirement including (Horn and Schunck 1981; Bruhn et al. 2005; Brox et al. 2004; Nir et al. 2008; Zimmer et al. 2009).

Equation (14) is then treated as a “calculus of variations” problem leading to the Euler-Lagrange equations:

$$\frac{\partial E_{\text{Global}}}{\partial u} - \frac{\partial}{\partial x} \frac{\partial E_{\text{Global}}}{\partial u_x} - \frac{\partial}{\partial y} \frac{\partial E_{\text{Global}}}{\partial u_y} = 0, \quad (15)$$

$$\frac{\partial E_{\text{Global}}}{\partial v} - \frac{\partial}{\partial x} \frac{\partial E_{\text{Global}}}{\partial v_x} - \frac{\partial}{\partial y} \frac{\partial E_{\text{Global}}}{\partial v_y} = 0. \quad (16)$$

Because they use the calculus of variations, such algorithms are generally referred to as *variational*. In the special case of the Horn-Schunck algorithm (Horn 1986), the Euler-Lagrange equations are linear in the unknown functions u and v . These equations are then parameterized with two unknown parameters per pixel and can be solved as a sparse linear system. A variety of options are possible, including the Jacobi method, the Gauss-Seidel method, Successive Over-Relaxation, and the Conjugate Gradient algorithm.

For more general energy functions, the Euler-Lagrange equations are nonlinear and are typically solved using an iterative method (analogous to gradient descent). For example, the flows can be parameterized by $u + du$ and $v + dv$ where u, v are treated as known (from the previous iteration or the initialization) and du, dv as unknowns. These expressions are substituted into the Euler-Lagrange equations, which are then linearized through the use of Taylor expansions. The resulting equations are linear in du and dv and solved using a sparse linear solver. The estimates of u and v are then updated appropriately and the next iteration applied.

One disadvantage of variational algorithms is that the discretization of the Euler-Lagrange equations is not always exact with respect to the original energy (Pock et al. 2007). Another extremal approach (Sun et al. 2008), closely related to the variational algorithms is to use:

$$\frac{\partial E_{\text{Global}}}{\partial \mathbf{f}} = 0 \quad (17)$$

rather than the Euler-Lagrange equations. Otherwise, the approach is similar. Equation (17) can be linearized and solved using a sparse linear system. The key difference between this approach and the variational one is just whether the parameterization of the flow functions into a set of flows per pixel occurs before or after the derivation of the extremal constraint equation ((17) or the Euler-Lagrange equations). One advantage of the early parameterization and the subsequent use of (17) is that it reduces the restrictions on the functional form of E_{Global} , important in learning-based approaches (Sun et al. 2008).

2.3.3 Other Continuous Algorithms

Another approach (Trobin et al. 2008; Wedel et al. 2008) is to decouple the data and prior terms through the introduction of two sets of flow parameters, say $(u_{\text{data}}, v_{\text{data}})$ for the data term and $(u_{\text{prior}}, v_{\text{prior}})$ for the prior:

$$E_{\text{Global}} = E_{\text{Data}}(u_{\text{data}}, v_{\text{data}}) + \lambda E_{\text{Prior}}(u_{\text{prior}}, v_{\text{prior}}) + \gamma (\|u_{\text{data}} - u_{\text{prior}}\|^2 + \|v_{\text{data}} - v_{\text{prior}}\|^2). \quad (18)$$

The final term in (18) encourages the two sets of flow parameters to be roughly the same. For a sufficiently large value

of γ the theoretical optimal solution will be unchanged and $(u_{\text{data}}, v_{\text{data}})$ will exactly equal $(u_{\text{prior}}, v_{\text{prior}})$. Practical optimization with too large a value of γ is problematic, however. In practice either a lower value is used or γ is steadily increased. The two sets of parameters allow the optimization to be broken into two steps. In the first step, the sum of the data term and the third term in (18) is optimized over the data flows $(u_{\text{data}}, v_{\text{data}})$ assuming the prior flows $(u_{\text{prior}}, v_{\text{prior}})$ are constant. In the second step, the sum of the prior term and the third term in (18) is optimized over prior flows $(u_{\text{prior}}, v_{\text{prior}})$ assuming the data flows $(u_{\text{data}}, v_{\text{data}})$ are constant. The result is two much simpler optimizations. The first optimization can be performed independently at each pixel. The second optimization is often simpler because it does not depend directly on the nonlinear data term (Trobin et al. 2008; Wedel et al. 2008).

Finally, in recent work, continuous convex optimization algorithms such as Linear Programming have also been used to compute optical flow (Seitz and Baker 2009).

2.3.4 Coarse-to-Fine and Other Heuristics

All of the above algorithms solve the problem as huge nonlinear optimizations. Even the Horn-Schunck algorithm, which results in linear Euler-Lagrange equations, is nonlinear through the linearization of the Brightness Constancy constraint to give the Optical Flow constraint. A variety of approaches have been used to improve the convergence rate and reduce the likelihood of falling into a local minimum.

One component in many algorithms is a coarse-to-fine strategy. The most common approach is to build image pyramids by repeated blurring and downsampling (Lucas and Kanade 1981; Glazer et al. 1983; Burt et al. 1983; Enkelman 1986; Anandan 1989; Black and Anandan 1996; Battiti et al. 1991; Bruhn et al. 2005). Optical flow is first computed on the top level (fewest pixels) and then upsampled and used to initialize the estimate at the next level. Computation at the higher levels in the pyramid involves far fewer unknowns and so is far faster. The initialization at each level from the previous level also means that far fewer iterations are required at each level. For this reason, pyramid algorithms tend to be significantly faster than a single solution at the bottom level. The images at the higher levels also contain fewer higher frequency components reducing the number of local minima in the data term. A related approach is to use a multigrid algorithm (Bruhn et al. 2006) where estimates of the flow are passed both up and down the hierarchy of approximations. A limitation of many coarse-to-fine algorithms, however, is the tendency to over-smooth fine structure and to fail to capture small fast-moving objects.

The main purpose of coarse-to-fine strategies is to deal with nonlinearities caused by the data term (and the subsequent difficulty in dealing with long-range motion). At the

coarsest pyramid level, the flow magnitude is likely to be small making the linearization of the brightness constancy assumption reasonable. Incremental warping of the flow between pyramid levels (Bergen et al. 1992) helps keep the flow update at any given level small (i.e., under one pixel). When combined with incremental warping and updating within a level, this method is effective for optimization with a linearized brightness constancy assumption.

Another common cause of nonlinearity is the use of a robust penalty function (see Sects. 2.1.2 and 2.2.2). A common approach to improve robustness in this case is Graduated Non-Convexity (GNC) (Blake and Zisserman 1987; Black and Anandan 1996). During GNC, the problem is first converted into a convex approximation that is more easily solved. The energy function is then made incrementally more non-convex and the solution is refined, until the original desired energy function is reached.

2.4 Discrete Optimization Algorithms

A number of recent approaches use discrete optimization algorithms, similar to those employed in stereo matching, such as graph cuts (Boykov et al. 2001) and belief propagation (Sun et al. 2003). Discrete optimization methods approximate the continuous space of solutions with a simplified problem. The hope is that this will enable a more thorough and complete search of the state space. The trade-off in moving from continuous to discrete optimization is one of search efficiency for fidelity. Note that, in contrast to discrete stereo optimization methods, the 2D flow field makes discrete optimization of optical flow significantly more challenging. Approximations are usually made, which can limit the power of the discrete algorithms to avoid local minima. The few methods proposed to date can be divided into two main approaches described below.

2.4.1 Fusion Approaches

Algorithms such as Jung et al. (2008), Lempitsky et al. (2008) and Trobin et al. (2008) assume that a number of candidate flow fields have been generated by running standard algorithms such as Lucas and Kanade (1981), and Horn and Schunck (1981), possibly multiple times with a number of different parameters. Computing the flow is then posed as choosing which of the set of possible candidates is best at each pixel. Fusion Flow (Lempitsky et al. 2008) uses a sequence of binary graph-cut optimizations to refine the current flow estimate by selectively replacing portions with one of the candidate solutions. Trobin et al. (2008) perform a similar sequence of fusion steps, at each step solving a continuous $[0, 1]$ optimization problem and then thresholding the results.

2.4.2 Dynamically Reparameterizing Sparse State-Spaces

Any fixed 2D discretization of the continuous space of 2D flow fields is likely to be a crude approximation to the continuous field. A number of algorithms take the approach of first approximating this state space sparsely (both spatially, and in terms of the possible flows at each pixel) and then refining the state space based on the result. An early use of this idea for flow estimation employed simulated annealing with a state space that adapted based on the local shape of the objective function (Black and Anandan 1991). More recently, Glocker et al. (2008) initially use a sparse sampling of possible motions on a coarse version of the problem. As the algorithm runs from coarse to fine, the spatial density of motion states (which are interpolated with a spline) and the density of possible flows at any given control point are chosen based on the uncertainty in the solution from the previous iteration. The algorithm of Lei and Yang (2009) also sparsely allocates states across space and for the possible flows at each spatial location. The spatial allocation uses a hierarchy of segmentations, with a single possible flow for each segment at each level. Within any level of the segmentation hierarchy, first a sparse sampling of the possible flows is used, followed by a denser sampling with a reduced range around the solution from the previous iteration. The algorithm in Cooke (2008) iteratively alternates between two steps. In the first step, all the states are allocated to the horizontal motion, which is estimated similarly to stereo, assuming the vertical motion is zero. In the second step, all the states are allocated to the vertical motion, treating the estimate of the horizontal motion from the previous iteration as constant.

2.4.3 Continuous Refinement

An optional step after a discrete algorithm is to use a continuous optimization to refine the results. Any of the approaches in Sect. 2.3 are possible.

2.5 Miscellaneous Issues

2.5.1 Learning

The design of a global energy function E_{Global} involves a variety of choices, each with a number of free parameters. Rather than manually making these decision and tuning parameters, learning algorithms have been used to choose the data and prior terms and optimize their parameters by maximizing performance on a set of training data (Roth and Black 2007; Sun et al. 2008; Li and Huttenlocher 2008).

2.5.2 Region-Based Techniques

If the image can be segmented into coherently moving regions, many of the methods above can be used to accu-

rately estimate the flow within the regions. Further, if the flow were accurately known, segmenting it into coherent regions would be feasible. One of the reasons optical flow has proven challenging to compute is that the flow and its segmentation must be computed together.

Several methods first segment the scene using non-motion cues and then estimate the flow in these regions (Black and Jepson 1996; Xu et al. 2008; Fuh and Maragos 1989). Within each image segment, Black and Jepson (1996) use a parametric model (e.g., affine) (Bergen et al. 1992), which simplifies the problem by reducing the number of parameters to be estimated. The flow is then refined as suggested above.

2.5.3 Layers

Motion transparency has been extensively studied and is not considered in detail here. Most methods have focused on the use of parametric models that estimate motion in layers (Jepson and Black 1993; Wang and Adelson 1993). The regularization of transparent motion in the framework of global energy minimization, however, has received little attention with the exception of Ju et al. (1996), Weiss (1997), and Shizawa and Mase (1991).

2.5.4 Sparse-to-Dense Approaches

The coarse-to-fine methods described above have difficulty dealing with long-range motion of small objects. In contrast, there exist many methods to accurately estimate sparse feature correspondences even when the motion is large. Such sparse matching method can be combined with the continuous energy minimization approaches in a variety of ways (Brox et al. 2009; Liu et al. 2008; Ren 2008; Xu et al. 2008).

2.5.5 Visibility and Occlusion

Occlusions and visibility changes can cause major problems for optical flow algorithms. The most common solution is to model such effects implicitly using a robust penalty function on both the data term and the prior term. Explicit occlusion estimation, for example through cross-checking flows computed forwards and backwards in time, is another approach that can be used to improve robustness to occlusions and visibility changes (Xu et al. 2008; Lei and Yang 2009).

2.6 Databases and Evaluations

Prior to our evaluation (Baker et al. 2007), there were three major attempts to quantitatively evaluate optical flow algo-

rithms, each proposing sequences with ground truth. The work of Barron et al. (1994) has been so influential that until recently, essentially all published methods compared with it. The synthetic sequences used there, however, are too simple to make meaningful comparisons between modern algorithms. Otte and Nagel (1994) introduced ground truth for a real scene consisting of polyhedral objects. While this provided real imagery, the images were extremely simple. More recently, McCane et al. (2001) provided ground truth for real polyhedral scenes as well as simple synthetic scenes. Most recently Liu et al. (2008) proposed a dataset of real imagery that uses hand segmentation and computed flow estimates within the segmented regions to generate the ground truth. While this has the advantage of using real imagery, the reliance on human judgement for segmentation, and on a particular optical flow algorithm for ground truth, may limit its applicability.

In this paper we go beyond these studies in several important ways. First, we provide ground-truth motion for much more complex real and synthetic scenes. Specifically, we include ground truth for scenes with nonrigid motion. Second, we also provide ground-truth motion boundaries and extend the evaluation methods to these areas where many flow algorithms fail. Finally, we provide a web-based interface, which facilitates the ongoing comparison of methods.

Our goal is to push the limits of current methods and, by exposing where and how they fail, focus attention on the hard problems. As described above, almost all flow algorithms have a specific data term, prior term, and optimization algorithm to compute the flow field. Regardless of the choices made, algorithms must somehow deal with all of the phenomena that make optical flow intrinsically ambiguous and difficult. These include: (1) the aperture problem and textureless regions, which highlight the fact that optical flow is inherently ill-posed, (2) camera noise, nonrigid motion, motion discontinuities, and occlusions, which make choosing appropriate penalty functions for both the data and prior terms important, (3) large motions and small objects which, often cause practical optimization algorithms to fall into local minima, and (4) mixed pixels, changes in illumination, non-Lambertian reflectance, and motion blur, which highlight overly simplified assumptions made by Brightness Constancy (or simple filter constancy). Our goal is to provide ground-truth data containing all of these components and to provide information about the location of motion boundaries and textureless regions. In this way, we hope to be able to evaluate which phenomena pose problems for which algorithms.

3 Database Design

Creating a ground-truth (GT) database for optical flow is difficult. For stereo, structured light (Scharstein and Szeliski

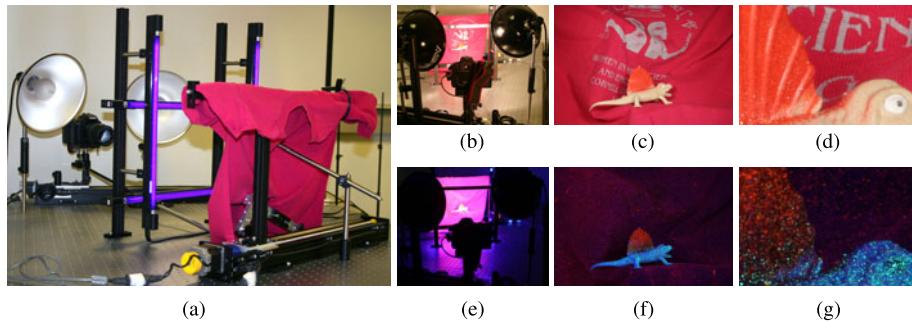


Fig. 1 (a) The setup for obtaining ground-truth flow using hidden fluorescent texture includes computer-controlled lighting to switch between the UV and visible lights. It also contains motion stages for both the camera and the scene. (b–d) The setup under the visible illumination. (e–g) The setup under the UV illumination. (c and f) Show the

high-resolution images taken by the digital camera. (d and g) Show a zoomed portion of (c) and (f). The high-frequency fluorescent texture in the images taken under UV light (g) allows accurate tracking, but is largely invisible in the low-resolution test images

2002) or range scanning (Seitz et al. 2006) can be used to obtain dense, pixel-accurate ground truth. For optical flow, the scene may be moving nonrigidly making such techniques inapplicable in general. Ideally we would like imagery collected in real-world scenarios with real cameras and substantial nonrigid motion. We would also like dense, subpixel-accurate ground truth. We are not aware of any technique that can simultaneously satisfy all of these goals.

Rather than collecting a single type of data (with its inherent limitations) we instead collected four different types of data, each satisfying a different subset of desirable properties. Having several different types of data has the benefit that the overall evaluation is less likely to be affected by any biases or inaccuracies in any of the data types. It is important to keep in mind that no ground-truth data is perfect. The term itself just means “measured on the ground” and any measurement process may introduce noise or bias. We believe that the combination of our four datasets is sufficient to allow a thorough evaluation of current optical flow algorithms. Moreover, the relative performance of algorithms on the different types of data is itself interesting and can provide insights for future algorithms (see Sect. 5.2.4).

Wherever possible, we collected eight frames with the ground-truth flow being defined between the middle pair. We collected color imagery, but also make grayscale imagery available for comparison with legacy implementations and existing approaches that only process grayscale. The dataset is divided into 12 training sequences with ground truth, which can be used for parameter estimation or learning, and 12 test sequences, where the ground truth is withheld. In this paper we only describe the test sequences. The datasets, instructions for evaluating results on the test set, and the performance of current algorithms are all available at <http://vision.middlebury.edu/flow/>. We describe each of the four types of data below.

3.1 Dense GT Using Hidden Fluorescent Texture

We have developed a technique for capturing imagery of nonrigid scenes with ground-truth optical flow. We build a scene that can be moved in very small steps by a computer-controlled motion stage. We apply a fine spatter pattern of fluorescent paint to all surfaces in the scene. The computer repeatedly takes a pair of high-resolution images both under ambient lighting and under UV lighting, and then moves the scene (and possibly the camera) by a small amount.

In our current setup, shown in Fig. 1(a), we use a Canon EOS 20D camera to take images of size 3504×2336 , and make sure that no scene point moves by more than 2 pixels from one captured frame to the next. We obtain our test sequence by downsampling every 40th image taken under visible light by a factor of six, yielding images of size 584×388 . Because we sample every 40th frame, the motion can be quite large (up to 12 pixels between frames in our evaluation data) even though the motion between each pair of captured frames is small and the frames are subsequently downsampled, i.e., after the downsampling, the motion between any pair of captured frames is at most $1/3$ of a pixel.

Since fluorescent paint is available in a variety of colors, the color of the objects in the scene can be closely matched. In addition, it is possible to apply a fine spatter pattern, where individual droplets are about the size of 1–2 pixels in the high-resolution images. This high-frequency texture is therefore far less perceptible in the low-resolution images, while the fluorescent paint is very visible in the high-resolution UV images in Fig. 1(g). Note that fluorescent paint absorbs UV light but emits light in the visible spectrum. Thus, the camera optics affect the hidden texture and the scene colors in exactly the same way, and the hidden texture remains perfectly aligned with the scene.

The ground-truth flow is computed by tracking small windows in the original sequence of high-resolution UV images. We use a sum-of-squared-difference (SSD) tracker

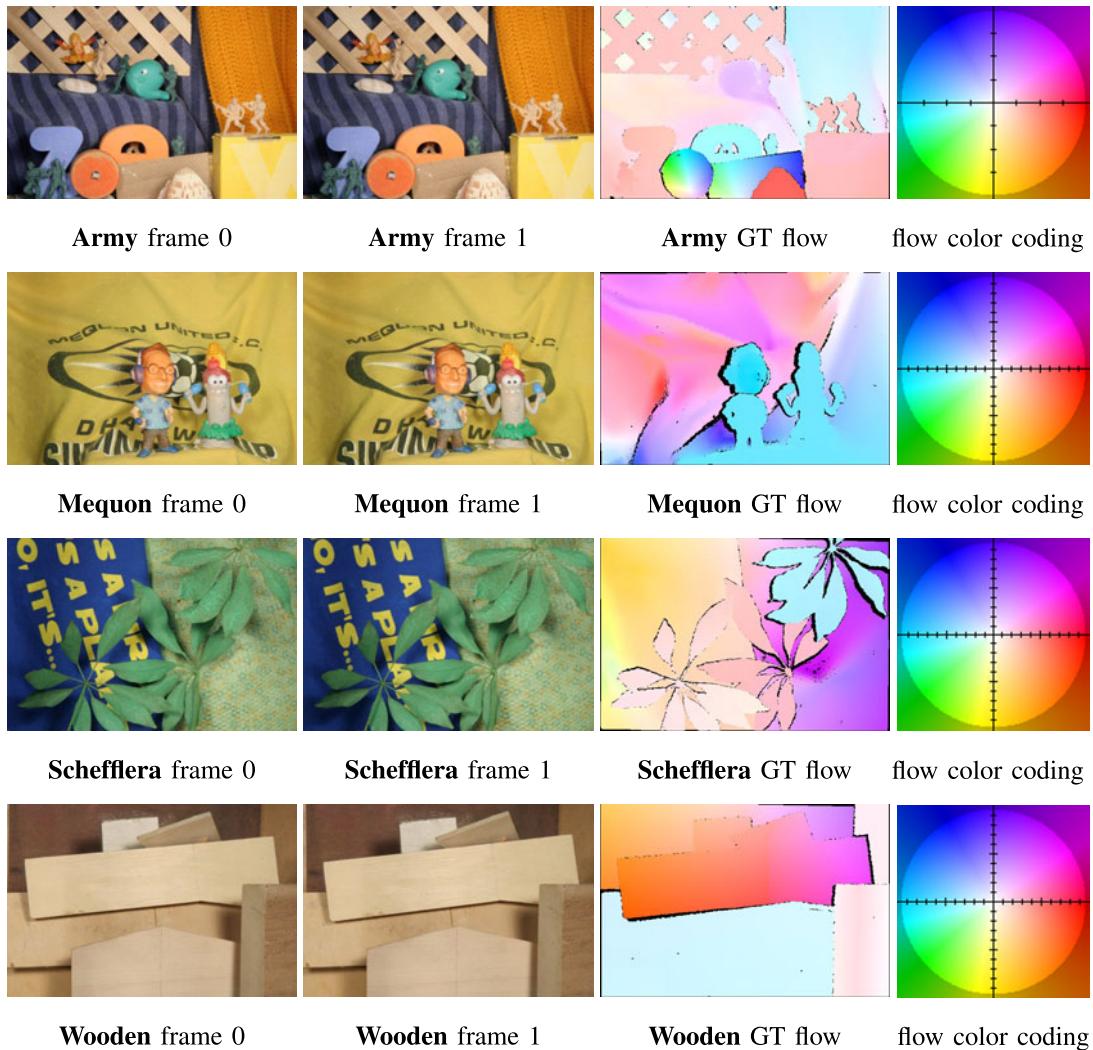


Fig. 2 Hidden Texture Data. *Army* contains several independently moving objects. *Mequon* contains nonrigid motion and textureless regions. *Schefflera* contains thin structures, shadows, and foreground/background transitions with little contrast. *Wooden* contains rigidly moving objects with little texture in the presence of shadows.

In the right-most column, we include a visualization of the color-coding of the optical flow. The “ticks” on the axes denote a flow unit of one pixel; note that the flow magnitudes are fairly low in *Army* (<4 pixels), but higher in the other three scenes (up to 10 pixels)

with a window size of 15×15 , corresponding to a window radius of less than 1.5 pixels in the downsampled images. We perform a local brute-force search, using each frame to initialize the next. We also crosscheck the results by tracking each pixel both forwards and backwards through the sequence and require perfect correspondence. The chances that this check would yield false positives after tracking for 40 frames are very low. Crosschecking identifies the occluded regions, whose motion we mark as “unknown.” After the initial integer-based motion tracking and crosschecking, we estimate the subpixel motion of each window using Lucas-Kanade (1981) with a precision of about 1/10 pixels (i.e., 1/60 pixels in the downsampled images). In order to downsample the motion field by a factor of 6, we find the modes among the 36 different motion vectors in each 6×6

window using sequential clustering. We assign the average motion of the dominant cluster as the motion estimate for the resulting pixel in the low-resolution motion field. The test images taken under visible light are downsampled using a binomial filter.

Using the combination of fluorescent paint, downsampling high-resolution images, and sequential tracking of small motions, we are able to obtain dense, subpixel accurate ground truth for a nonrigid scene.

We include four sequences in the evaluation set (Fig. 2). *Army* contains several independently moving objects. *Mequon* contains nonrigid motion and large areas with little texture. *Schefflera* contains thin structures, shadows, and foreground/background transitions with little contrast. *Wooden* contains rigidly moving objects with little texture

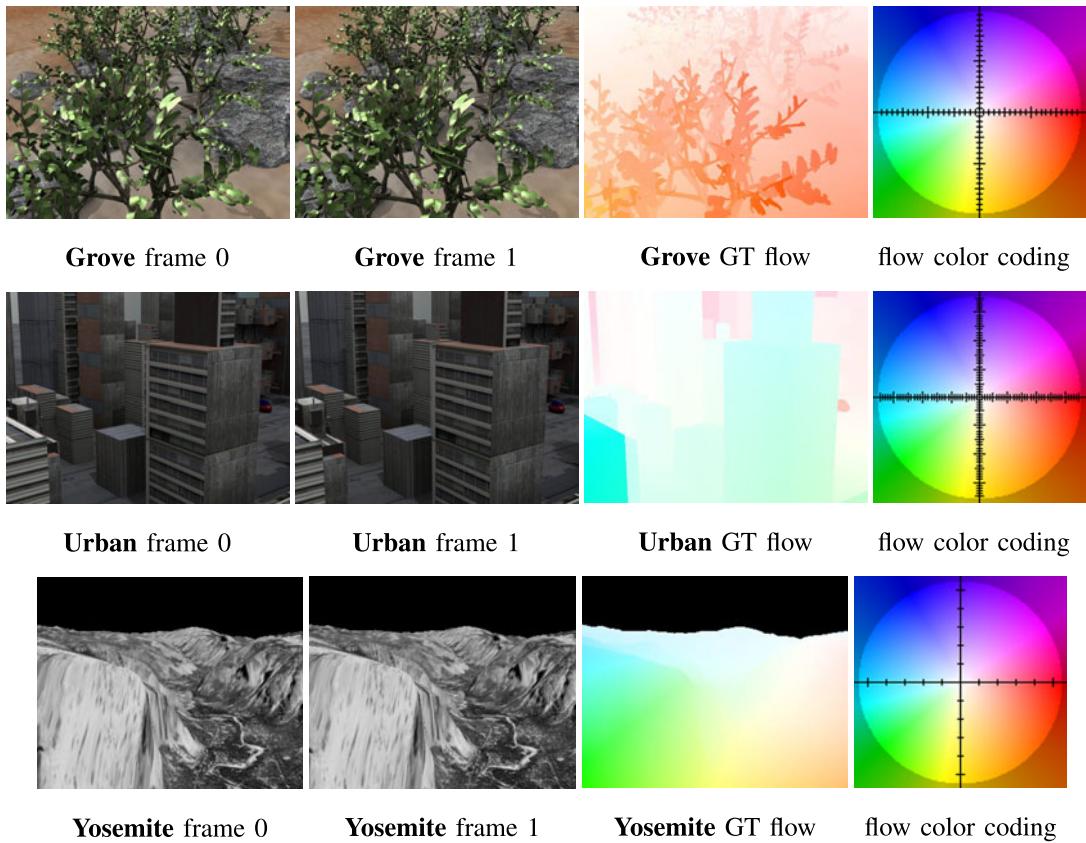


Fig. 3 Synthetic Data. *Grove* contains a close up of a tree with thin structures, very complex motion discontinuities, and a large motion range (up to 20 pixels). *Urban* contains large motion discontinuities

and an even larger motion range (up to 35 pixels). *Yosemite* is included in our evaluation to allow comparison with algorithms published prior to our study

in the presence of shadows. The maximum motion in *Army* is approximately 4 pixels. The maximum motion in the other three sequences is about 10 pixels. All sequences are significantly more difficult than the *Yosemite* sequence due to the larger motion ranges, the non-rigid motion, various photometric effects such as shadows and specularities, and the detailed geometric structure.

The main benefit of this dataset is that it contains ground truth on imagery captured with a real camera. Hence, it contains real photometric effects, natural textural properties, etc. The main limitations of this dataset are that the scenes are laboratory scenes, not real-world scenes. There is also no motion blur due to the stop motion method of capture.

One drawback of this data is that the ground truth it is not available in areas where cross-checking failed, in particular, in regions occluded in one image. Even though the ground truth is reasonably accurate (on the order of 1/60th of a pixel), the process is not perfect; significant errors however, are limited to a small fraction of the pixels. The same can be said for any real data where the ground truth is measured, including, for example, in the Middlebury stereo dataset (Scharstein and Szeliski 2002). The ground-truth measuring

technique may always be prone to errors and biases. Consequently, the following section describes realistic synthetic data where the ground truth is guaranteed to be perfect.

3.2 Realistic Synthetic Imagery

Synthetic scenes generated using computer graphics are often indistinguishable from real ones. For the study of optical flow, synthetic data offers a number of benefits. In particular, it gives full control over the rendering process including material properties of the objects, while providing precise ground-truth motion and object boundaries.

To go beyond previous synthetic ground truth (e.g., the *Yosemite* sequence), we generated two types of fairly complex synthetic outdoor scenes. The first is a set of “natural” scenes (Fig. 3 top) containing significant complex occlusion. These scenes consist of a random number of procedurally generated rocks and trees with randomly chosen ground texture and surface displacement. Additionally, the tree bark has significant 3D texture. The trees have a small amount of independent movement to mimic motion due to wind. The camera motions include camera rotation and 3D translation. A second set of “urban” scenes (Fig. 3 middle) con-

tain buildings generated with a random shape grammar. The buildings have randomly selected scanned textures; there are also a few independently moving “cars.”

These scenes were generated using the 3Delight Renderman-compliant renderer (DNA Research 2008) at a resolution of 640×480 pixels using linear gamma. The images are antialiased, mimicking the effect of sensors with finite area. Frames in these synthetic sequences were generated without motion blur. There are cast shadows, some of which are non-stationary due to the independent motion of the trees and cars. The surfaces are mostly diffuse, but the leaves on the trees have a slight specular component, and the cars are strongly specular. A minority of the surfaces in the urban scenes have a small (5%) reflective component, meaning that the reflection of other objects is faintly visible in these surfaces.

The rendered scenes use the ambient occlusion approximation to global illumination (Landis 2002). This approximation separates illumination into the sum of direct and multiple-bounce components, and then assumes that the multiple-bounce illumination is sufficiently omnidirectional that it can be approximated at each point by a product of the incoming ambient light and a precomputed factor measuring the proportion of rays that are not blocked by other nearby surfaces.

The ground truth was computed using a custom shader that projects the 3D motion of the scene corresponding to a particular image onto the 2D image plane. Since individual pixels can potentially represent more than one object, simply point-sampling the flow at the center of each pixel could result in a flow vector that does not reflect the dominant motion under the pixel. On the other hand, applying antialiasing to the flow would result in an averaged flow vector at each pixel that does reflect the true motion of any object within that pixel. Instead, we clustered the flow vectors within each pixel and selected a flow vector from the dominant cluster. The flow fields are initially generated at $3 \times$ resolution, resulting in nine candidate flow vectors for each pixel. These motion vectors are grouped into two clusters using k -means. The k -means procedure is initialized with the vectors closest and furthest from the pixel’s average flow as measured using the flow vector end points. The flow vector closest to the mean of the dominant cluster is then chosen to represent the flow for that pixel. The images were also generated at $3 \times$ resolution and downsampled using a bicubic filter.

We selected three synthetic sequences to include in the evaluation set (Fig. 3). *Grove* contains a close-up view of a tree, with a substantial parallax and motion discontinuities. *Urban* contains images of a city, with substantial motion discontinuities, a large motion range, and an independently moving object. We also include the *Yosemite* sequence to allow some comparison with algorithms published prior to the release of our data.

3.3 Imagery for Frame Interpolation

In a wide class of applications such as video re-timing, novel view generation, and motion-compensated compression, what is important is not how well the flow field matches the ground-truth motion, but how well intermediate frames can be predicted using the flow. To allow for measures that predict performance on such tasks, we collected a variety of data suitable for frame interpolation. The relative performance of algorithms with respect to frame interpolation and ground-truth motion estimation is interesting in its own right.

3.3.1 Frame Interpolation Datasets

We used a PointGrey Dragonfly Express camera to capture the data, acquiring 60 frames per second. We provide every other frame to the optical flow algorithms and retain the intermediate images as frame-interpolation ground truth. This temporal subsampling means that the input to the flow algorithms is captured at 30 Hz while enabling generation of a $2 \times$ slow-motion sequence.

We include four such sequences in the evaluation set (Fig. 4). The first two (*Backyard* and *Basketball*) include people, a common focus of many applications, but a subject matter absent from previous evaluations. *Backyard* is captured outdoors with a short shutter (6 ms) and has little motion blur. *Basketball* is captured indoors with a longer shutter (16 ms) and so has more motion blur. The third sequence, *Dumptruck*, is an urban scene containing several independently moving vehicles, and has substantial specularities and saturation (2 ms shutter). The final sequence, *Evergreen*, includes highly textured vegetation with complex motion discontinuities (6 ms shutter).

The main benefit of the interpolation dataset is that the scenes are real world scenes, captured with a real camera and containing real sources of noise. The ground truth is not a flow field, however, but an intermediate image frame. Hence, the definition of flow being used is the *apparent motion*, not the 2D projection of the motion field.

3.3.2 Frame Interpolation Algorithm

Note that the evaluation of accuracy depends on the interpolation algorithm used to construct the intermediate frame. By default, we generate the intermediate frames from the flow fields uploaded to the website using our baseline interpolation algorithm. Researchers can also upload their own interpolation results in case they want to use a more sophisticated algorithm.

Our algorithm takes a single flow field \mathbf{u}_0 from image I_0 to I_1 and constructs an interpolated frame I_t at time $t \in (0, 1)$. We do, however, use both frames to generate the



Fig. 4 High-Speed Data for Interpolation. We collected four sequences using a PointGrey Dragonfly Express running at 60 Hz. We provide every other image to the algorithms and retain the intermediate frame as interpolation ground truth. The first two sequences (*Backyard*

and *Basketball*) include people, a common focus of many applications. *Dumptruck* contains several independently moving vehicles, and has substantial specularities and saturation. *Evergreen* includes highly textured vegetation with complex discontinuities

actual intensity values. In all the experiments in this paper $t = 0.5$. Our algorithm is closely related to previous algorithms for depth-based frame interpolation (Shade et al. 1998; Zitnick et al. 2004):

- (1) Forward-warp the flow \mathbf{u}_0 to time t to give \mathbf{u}_t where:

$$\mathbf{u}_t(\text{round}(\mathbf{x} + t\mathbf{u}_0(\mathbf{x}))) = \mathbf{u}_0(\mathbf{x}). \quad (19)$$

In order to avoid sampling gaps, we splat the flow vectors with a splatting radius of ± 0.5 pixels (Levoy 1988) (i.e., each flow vector is followed to a real-valued location in the destination image, and the flow is written into all pixels within a distance of 0.5 of that location). In cases where multiple flow vectors map to the same location, we attempt to resolve the ordering indepen-

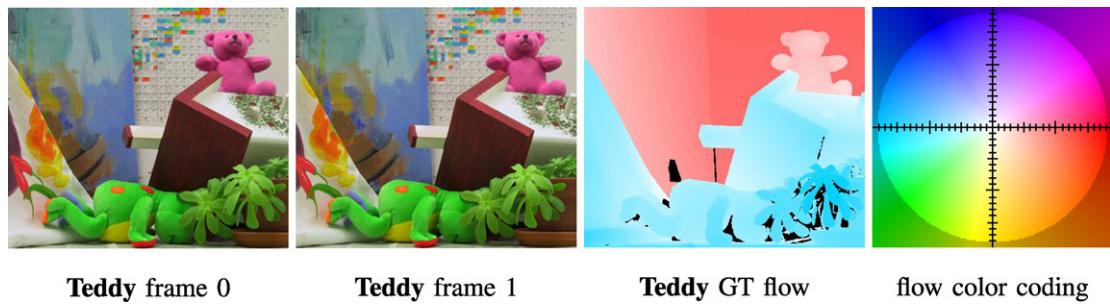


Fig. 5 Stereo Data. We cropped the stereo dataset *Teddy* (Scharstein and Szeliski 2003) to convert the asymmetric stereo disparity range into a roughly symmetric flow field. This dataset includes complex

geometry as well as significant occlusions and motion discontinuities. One reason for including this dataset is to allow comparison with state-of-the-art stereo algorithms

dently for each pixel by checking photoconsistency; i.e., we retain the flow $\mathbf{u}_0(\mathbf{x})$ with the lowest color difference $|I_0(\mathbf{x}) - I_1(\mathbf{x} + \mathbf{u}_0(\mathbf{x}))|$.

- (2) Fill any holes in \mathbf{u}_t using a simple outside-in strategy.
- (3) Estimate occlusion masks $O_0(\mathbf{x})$ and $O_1(\mathbf{x})$, where $O_i(\mathbf{x}) = 1$ means pixel \mathbf{x} in image I_i is not visible in the respective other image. To compute $O_0(\mathbf{x})$ and $O_1(\mathbf{x})$, we first forward-warp the flow $\mathbf{u}_0(\mathbf{x})$ to time $t = 1$ using the same approach as in Step 1 to give $\mathbf{u}_1(\mathbf{x})$. Any pixel \mathbf{x} in $\mathbf{u}_1(\mathbf{x})$ that is not targeted by this splatting has no corresponding pixel in I_0 and thus we set $O_1(\mathbf{x}) = 1$ for all such pixels. (See Herbst et al. 2009 for a bidirectional algorithm that performs this reasoning at time t .) In order to compute $O_0(\mathbf{x})$, we cross-check the flow vectors, setting $O_0(\mathbf{x}) = 1$ if

$$|\mathbf{u}_0(\mathbf{x}) - \mathbf{u}_1(\mathbf{x} + \mathbf{u}_0(\mathbf{x}))| > 0.5. \quad (20)$$

- (4) Compute the colors of the interpolated pixels, taking occlusions into consideration. Let $\mathbf{x}_0 = \mathbf{x} - t\mathbf{u}_t(\mathbf{x})$ and $\mathbf{x}_1 = \mathbf{x} + (1 - t)\mathbf{u}_t(\mathbf{x})$ denote the locations of the two “source” pixels in the two images. If both pixels are visible, i.e., $O_0(\mathbf{x}_0) = 0$ and $O_1(\mathbf{x}_1) = 0$, blend the two images (Beier and Neely 1992):

$$I_t(\mathbf{x}) = (1 - t)I_0(\mathbf{x}_0) + tI_1(\mathbf{x}_1). \quad (21)$$

Otherwise, only sample the non-occluded image, i.e., set $I_t(\mathbf{x}) = I_0(\mathbf{x}_0)$ if $O_1(\mathbf{x}_1) = 1$ and vice versa. In order to avoid artifacts near object boundaries, we dilate the occlusion masks O_0 , O_1 by a small radius before this operation. We use bilinear interpolation to sample the images.

This algorithm, while reasonable, is only meant to serve as starting point. One area for future research is to develop better frame interpolation algorithms. We hope that our database will be used both by researchers working on opti-

cal flow and on frame interpolation (Mahajan et al. 2009; Herbst et al. 2009).

3.4 Modified Stereo Data for Rigid Scenes

Our final type of data consists of modified stereo data. Specifically we include the *Teddy* dataset in the evaluation set, the ground truth for which was obtained using structured lighting (Scharstein and Szeliski 2003) (Fig. 5). Stereo datasets typically have an asymmetric disparity range $[0, d_{\max}]$, which is appropriate for stereo, but not for optical flow. We crop different subregions of the images, thereby introducing a spatial shift, to convert this disparity range to $[-d_{\max}/2, d_{\max}/2]$.

A key benefit of the modified stereo dataset, like the hidden fluorescent texture dataset, is that it contains ground-truth flow fields on imagery captured with a real camera. An additional benefit is that it allows a comparison between state-of-the-art stereo algorithms and optical flow algorithms (see Sect. 5.6). Shifting the disparity range does not affect the performance of stereo algorithms as long as they are given the new search range. Although optical flow is a more under-constrained problem, the relative performance of algorithms may lead to algorithmic insights.

One concern with the modified stereo dataset is that algorithms may take advantage of the knowledge that the motions are all horizontal. Indeed a number recent algorithms have considered rigidity priors (Wedel et al. 2008, 2009). However, these algorithms must also perform well on the other types of data and any over-fitting to the rigid data should be visible by comparing results across the 12 images in the evaluation set. Another concern would be that the ground truth is only accurate to 0.25 pixels. (The original stereo data comes with pixel-accurate ground truth but is four times higher resolution—Scharstein and Szeliski 2003.) The most appropriate performance statistics for this data, therefore, are the robustness statistics used in the Middlebury stereo dataset (Scharstein and Szeliski 2002) (Sect. 4.2).

4 Evaluation Methodology

We refine and extend the evaluation methodology of Barron et al. (1994) in terms of: (1) the performance measures used, (2) the statistics computed, and (3) the sub-regions of the images considered.

4.1 Performance Measures

The most commonly used measure of performance for optical flow is the angular error (AE). The AE between a flow vector (u, v) and the ground-truth flow $(u_{\text{GT}}, v_{\text{GT}})$ is the angle in 3D space between $(u, v, 1.0)$ and $(u_{\text{GT}}, v_{\text{GT}}, 1.0)$. The AE can be computed by taking the dot product of the vectors, dividing by the product of their lengths, and then taking the inverse cosine:

$$\text{AE} = \cos^{-1} \left(\frac{1.0 + u \times u_{\text{GT}} + v \times v_{\text{GT}}}{\sqrt{1.0 + u^2 + v^2} \sqrt{1.0 + u_{\text{GT}}^2 + v_{\text{GT}}^2}} \right). \quad (22)$$

The popularity of this measure is based on the seminal survey by Barron et al. (1994), although the measure itself dates to prior work by Fleet and Jepson (1990). The goal of the AE is to provide a *relative* measure of performance that avoids the “divide by zero” problem for zero flows. Errors in large flows are penalized less in AE than errors in small flows.

Although the AE is prevalent, it is unclear why errors in a region of smooth non-zero motion should be penalized less than errors in regions of zero motion. The AE also contains an arbitrary scaling constant (1.0) to convert the units from pixels to degrees. Hence, we also compute an *absolute* error, the error in flow endpoint (EE) used in Otte and Nagel (1994) defined by:

$$\text{EE} = \sqrt{(u - u_{\text{GT}})^2 + (v - v_{\text{GT}})^2}. \quad (23)$$

Although the use of AE is common, the EE measure is probably more appropriate for most applications (see Sect. 5.2.1). We report both.

For image interpolation, we define the interpolation error (IE) to be the root-mean-square (RMS) difference between the ground-truth image and the estimated interpolated image

$$\text{IE} = \left[\frac{1}{N} \sum_{(x,y)} (I(x, y) - I_{\text{GT}}(x, y))^2 \right]^{\frac{1}{2}}, \quad (24)$$

where N is the number of pixels. For color images, we take the L2 norm of the vector of RGB color differences.

We also compute a second measure of interpolation performance, a gradient-normalized RMS error inspired by

Szeliski (1999). The normalized interpolation error (NE) between an interpolated image $I(x, y)$ and a ground-truth image $I_{\text{GT}}(x, y)$ is given by:

$$\text{NE} = \left[\frac{1}{N} \sum_{(x,y)} \frac{(I(x, y) - I_{\text{GT}}(x, y))^2}{\|\nabla I_{\text{GT}}(x, y)\|^2 + \epsilon} \right]^{\frac{1}{2}}. \quad (25)$$

In our experiments the arbitrary scaling constant is set to be $\epsilon = 1.0$ (graylevels per pixel squared). Again, for color images, we take the L2 norm of the vector of RGB color differences and compute the gradient of each color band separately.

Naturally, an interpolation algorithm is required to generate the interpolated image from the optical flow field. In this paper, we use the baseline algorithm outlined in Sect. 3.3.2.

4.2 Statistics

Although the full histograms are available in a technical report, Barron et al. (1994) only reports averages (AV) and standard deviations (SD). This has led most subsequent researchers to only report these statistics. We also compute the robustness statistics used in the Middlebury stereo dataset (Scharstein and Szeliski 2002). In particular RX denotes the percentage of pixels that have an error measure above X. For the angle error (AE) we compute R2.5, R5.0, and R10.0 (degrees); for the endpoint error (EE) we compute R0.5, R1.0, and R2.0 (pixels); for the interpolation error (IE) we compute R2.5, R5.0, and R10.0 (graylevels); and for the normalized interpolation error (NE) we compute R0.5, R1.0, and R2.0 (no units). We also compute robust accuracy measures similar to those in Seitz et al. (2006): AX denotes the accuracy of the error measure at the Xth percentile, after sorting the errors from low to high. For the flow errors (AE and EE), we compute A50, A75, and A95. For the interpolation errors (IE and NE), we compute A90, A95, and A99.

4.3 Region Masks

It is easier to compute flow in some parts of an image than in others. For example, computing flow around motion discontinuities is hard. Computing motion in textureless regions is also hard, although interpolating in those regions should be easier. Computing statistics over such regions may highlight areas where existing algorithms are failing and spur further research in these cases. We follow the procedure in Scharstein and Szeliski (2002) and compute the error measure statistics over three types of region masks: everywhere (*All*), around motion discontinuities (*Disc*), and in textureless regions (*Untext*). We illustrate the masks for the *Scheflera* dataset in Fig. 6.

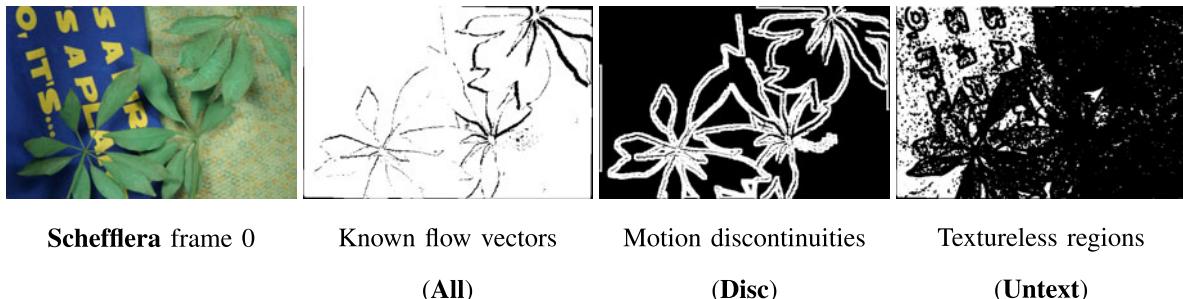


Fig. 6 Region masks for *Schefflera*. Statistics are computed over the white pixels. *All* includes all the pixels where the ground-truth flow can be reliably determined. The *Disc* mask is computed by taking the gradient of the ground-truth flow (or pixel differencing if the ground-

Motion discontinuities
(**Disc**) Textureless regions
(**Untext**)

truth flow is unavailable), thresholding and dilating. The *Untext* regions are computed by taking the gradient of the image, thresholding and dilating

The *All* masks for flow estimation include all the pixels where the ground-truth flow could be reliably determined. For the new synthetic sequences, this means all of the pixels. For *Yosemite*, the sky is excluded. For the hidden fluorescent texture data, pixels where cross-checking failed are excluded. Most of these pixels are around the boundary of objects, and around the boundary of the image where the pixel flows outside the second image. Similarly, for the stereo sequences, pixels where cross-checking failed are excluded (Scharstein and Szeliski 2003). Most of these pixels are pixels that are occluded in one of the images. The *All* masks for the interpolation metrics include all of the pixels. Note that in some cases (particularly the synthetic data), the *All* masks include pixels that are visible in first image but are occluded or outside the second image. We did not remove these pixels because we believe algorithms should be able to extrapolate into these regions.

The *Disc* mask is computed by taking the gradient of the ground-truth flow field, thresholding the magnitude, and then dilating the resulting mask with a 9×9 box. If the ground-truth flow is not available, we use frame differencing to get an estimate of fast-moving regions instead. The *Untext* regions are computed by taking the gradient of the image, thresholding the magnitude, and dilating with a 3×3 box. The pixels excluded from the *All* masks are also excluded from both *Disc* and *Untext* masks.

5 Experimental Results

We now discuss our empirical findings. We start in Sect. 5.1 by outlining the evolution of our online evaluation since the publication of our preliminary paper (Baker et al. 2007). In Sect. 5.2, we analyze the flow errors. In particular, we investigate the correlation between the various metrics, statistics, region masks, and datasets. In Sect. 5.3, we analyze the interpolation errors and in Sect. 5.4, we compare the interpolation error results with the flow error results. Finally,

in Sect. 5.5, we compare the algorithms that have reported results using our evaluation in terms of which components of our taxonomy in Sect. 2 they use.

5.1 Online Evaluation

Our online evaluation at <http://vision.middlebury.edu/flow/> provides a snapshot of the state-of-the-art in optical flow. Seeded with the handful of methods that we implemented as part of our preliminary paper (Baker et al. 2007), the evaluation has quickly grown. At the time of writing (December 2009), the evaluation contains results for 24 published methods and several unpublished ones. In this paper, we restrict attention to the published algorithms. Four of these methods were contributed by us (our implementations of Horn and Schunck 1981, Lucas-Kanade 1981, Combined Local-Global—Bruhn et al. 2005, and Black and Anandan 1996). Results for the 20 other methods were submitted by their authors. Of these new algorithms, two were published before 2007, 11 were published in 2008, and 7 were published in 2009.

On the evaluation website, we provide tables comparing the performance of the algorithms for each of the four error measures, i.e., endpoint error (EE), angular error (AE), interpolation error (IE), and normalized interpolation error (NE), on a set of 8 test sequences. For EE and AE, which measure flow accuracy, we use the 8 sequences for which we have ground-truth flow: *Army*, *Mequon*, *Schefflera*, *Wooden*, *Grove*, *Urban*, *Yosemite*, and *Teddy*. For IE and NE, which measure interpolation accuracy, we use only four of the above datasets (*Mequon*, *Schefflera*, *Urban*, and *Teddy*) and replace the other four with the high-speed datasets *Backyard*, *Basketball*, *Dumptruck*, and *Evergreen*. For each measure, we include a separate page for each of the eight statistics in Sect. 4.2. Figure 7 shows a screenshot of the first of these 32 pages, the average endpoint error (Avg. EE). For each measure and statistic, we evaluate all methods on the set of eight test images with three different regions masks

Optical flow evaluation results			Statistics:		Average	SD	R0.5	R1.0	R2.0	A50	A75	A95
			Error type:		endpoint	angle	interpolation	normalized interpolation				
Average endpoint error	Army (Hidden texture)	Mequon (Hidden texture)	Schefflera (Hidden texture)	Wooden (Hidden texture)	Grove (Synthetic)	Urban (Synthetic)	Yosemite (Synthetic)	Teddy (Stereoscopic)				
avg. rank	GT im0 im1	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
Adaptive [20]	4.4	0.09 1.26 1.06 1	0.23 0.78 0.05 1	0.54 0.11 0.05 1	0.21 0.18 0.10 1	0.88 1.25 0.73 1	0.50 1.28 0.31 1	0.14 10.16 12.22 10	0.65 3	1.37 3	0.79 4	
Complementary OF [21]	5.7	0.11 0.28 0.10 9	0.18 1.63 0.12 1	0.21 0.13 0.75 1	0.18 0.19 0.75 1	0.97 10.13 1.6 1.00 11	0.78 20.17.37 0.87 14	0.11 4 0.12 2 0.22 10	0.68 4	1.48 4	0.95 8	
Aniso-Huber-L1 [22]	5.8	0.10 0.28 0.08 3	0.31 0.88 0.28 12	0.10 1.17 0.29 12	0.20 4 0.92 0.13 5	0.84 2 1.20 2 0.70 2	0.39 1 1.23 1 0.28 1	0.17 15 0.15 9 0.27 16	0.64 2	1.36 2	0.79 4	
DPOF [18]	6.1	0.13 12.35 12.09 4	0.25 6 0.79 5	0.17 0.24 0.19 1	0.21 3 0.62 1 0.15 11	0.74 1 1.09 1 0.49 1	0.66 7 1.80 10 0.63 8	0.19 17 0.17 14 0.35 20	0.50 1	1.08 1	0.55 1	
TV-L1-improved [17]	7.2	0.09 1.26 1.07 2	0.20 3 0.71 3 0.16 2	0.53 7 1.18 3 0.22 8	0.21 7 1.24 11 0.11 2	0.90 4 1.31 6 0.72 3	1.51 14 1.93 11 0.84 11	0.18 16 0.17 14 0.31 17	0.73 8	1.62 9	0.87 7	
CBF [12]	7.8	0.10 0.28 0.09 4	0.34 12 0.80 6 0.37 13	0.43 5 0.95 5 0.26 21	0.21 7 1.14 8 0.13 5	0.90 4 1.27 4 0.82 7	0.41 2 1.23 1 0.30 2	0.23 22 0.19 20 0.39 21	0.76 9	1.56 6	1.02 9	
Brox et al. [5]	8.4	0.11 5 0.32 8 0.11 12	0.27 9 0.93 10 0.22 9	0.39 4 0.94 4 0.24 7	0.24 9 1.25 12 0.13 5	1.10 13 1.39 12 1.43 17	0.88 6 1.77 8 0.55 7	0.10 2 0.13 4 0.11 1	0.91 11	1.83 12	1.13 12	
Rannacher [23]	8.5	0.11 5 0.31 6 0.09 4	0.25 6 0.84 7 0.21 8	0.57 12 1.27 15 0.26 8	0.24 9 1.32 14 0.13 5	0.91 7 1.33 8 0.72 3	1.49 13 1.95 13 0.78 9	0.15 12 0.14 7 0.26 13	0.69 6	1.58 8	0.86 6	
F-TV-L1 [15]	8.8	0.14 13 0.35 12 0.14 15	0.34 12 0.98 12 0.26 11	0.59 14 1.19 10 0.26 8	0.27 13 1.36 15 0.16 12	0.90 4 1.30 5 0.76 6	0.54 4 1.62 6 0.36 4	0.13 6 0.15 9 0.20 9	0.68 4	1.56 6	0.86 2	
Second-order prior [8]	9.0	0.11 5 0.31 6 0.09 4	0.26 8 0.93 10 0.20 7	0.57 12 1.25 14 0.26 8	0.20 4 1.04 6 0.12 3	0.94 8 1.34 9 0.83 9	0.61 6 1.93 11 0.47 8	0.20 18 0.16 12 0.34 19	0.77 10	1.64 10	1.07 10	
Fusion [6]	9.4	0.11 5 0.34 10 0.10 9	0.19 2 0.69 2 0.16 2	0.28 6 0.66 2 0.23 6	0.20 4 1.19 8 0.14 9	1.07 11 1.42 13 1.22 13	1.35 10 1.49 8 0.86 13	0.20 18 0.20 21 0.26 13	1.07 14	2.07 16	1.39 16	
Dynamic MRF [7]	11.1	0.12 11 0.34 10 0.11 12	0.22 4 0.89 9 0.16 2	0.44 6 1.13 7 0.20 2	0.24 9 1.29 13 0.14 9	1.11 14 1.52 17 1.13 12	1.54 15 2.37 20 0.93 18	0.13 6 0.12 2 0.31 17	1.27 18	2.33 20	1.66 17	
SegOF [10]	11.7	0.15 14 0.36 14 0.10 9	0.57 15 1.16 15 0.59 19	0.68 15 1.24 12 0.64 14	0.32 15 0.86 2 0.26 15	1.18 17 1.50 16 1.47 18	1.63 18 2.09 14 0.98 16	0.08 1 0.13 4 0.12 2	0.70 7	1.50 5	0.89 3	
Learning Flow [11]	13.3	0.11 5 0.32 8 0.09 4	0.29 10 0.99 13 0.23 10	0.55 9 1.24 12 0.29 10	0.36 16 1.56 17 0.25 14	1.25 19 1.61 21 1.41 16	1.55 17 2.32 19 0.85 12	0.14 10 0.18 18 0.24 12	1.09 15	2.09 18	1.27 13	
Filter Flow [19]	14.3	0.17 16 0.39 16 0.13 14	0.43 14 1.09 14 0.38 14	0.75 16 1.34 16 0.78 19	0.70 19 1.54 16 0.68 19	1.33 16 1.38 11 1.51 19	0.57 6 1.32 4 0.44 6	0.22 20 0.23 23 0.26 13	0.96 12	1.66 11	1.12 11	
GraphCuts [14]	14.5	0.16 15 0.38 15 0.14 15	0.59 18 1.36 19 0.46 15	0.56 10 1.07 6 0.64 14	0.26 12 1.14 8 0.17 13	0.96 9 1.35 10 0.84 10	2.25 23 1.79 9 1.22 21	0.22 20 0.17 14 0.43 22	1.22 17	2.05 15	1.78 19	
Black & Anandan [4]	15.0	0.18 17 0.42 17 0.19 18	0.58 17 1.31 17 0.50 16	0.95 19 1.58 19 0.70 16	0.49 17 1.59 18 0.45 17	1.08 12 1.42 13 1.22 13	1.43 11 2.28 17 0.83 10	0.15 12 0.17 14 0.17 6	1.11 16	1.98 14	1.30 14	
SPSA-learn [13]	15.7	0.18 17 0.45 18 0.17 17	0.57 15 1.32 18 0.51 17	0.84 17 1.50 17 0.72 17	0.52 18 1.64 19 0.49 18	1.12 15 1.42 13 1.39 19	1.75 19 2.14 18 1.06 20	0.13 6 0.13 4 0.19 7	1.32 19	2.08 17	1.73 18	
GroupFlow [9]	15.8	0.21 19 0.51 19 0.21 19	0.79 21 1.69 21 0.72 21	0.86 18 1.64 19 0.74 18	0.30 14 1.07 2 0.26 15	1.29 22 1.81 22 0.82 7	1.94 21 2.30 18 1.36 22	0.11 4 0.14 7 0.19 7	1.06 13	1.96 13	1.35 15	
2D-CLG [1]	17.4	0.28 21 0.62 22 0.21 19	0.67 20 1.21 16 0.70 20	1.12 21 1.80 21 0.98 22	1.07 22 2.06 21 1.12 22	1.23 18 1.52 17 1.62 22	1.54 18 2.15 16 0.96 16	0.10 2 0.11 1 0.16 4	1.38 20	2.26 19	1.83 20	
Horn & Schunck [3]	18.6	0.22 20 0.55 20 0.22 21	0.61 19 1.53 20 0.52 18	1.01 20 1.73 20 0.80 20	0.78 20 2.02 20 0.77 20	1.26 20 1.58 19 1.55 20	1.43 11 2.59 22 1.00 18	0.16 14 0.18 18 0.15 3	1.51 21	2.50 21	1.88 21	
Ti-DOFE [24]	19.6	0.38 23 0.64 23 0.47 23	1.16 22 1.72 22 1.26 22	1.39 23 2.06 24 1.17 23	1.29 23 2.21 23 1.41 23	1.27 21 1.61 20 1.57 21	1.28 9 2.57 21 1.01 19	0.13 6 0.15 9 0.16 4	1.87 22	2.71 22	2.53 22	
FOLKI [16]	22.6	0.29 22 0.73 24 0.33 22	1.52 23 1.96 24 1.80 23	1.23 22 2.04 22 0.95 21	0.98 21 2.20 22 1.08 21	1.53 23 1.85 23 2.07 23	2.14 22 3.23 24 1.60 23	0.26 23 0.21 22 0.68 23	2.67 23	3.27 23	4.32 23	
Pyramid LK [2]	23.7	0.39 24 0.61 21 0.61 24	1.67 24 1.78 23 2.00 24	1.50 24 1.97 22 1.38 24	1.57 24 2.39 24 1.78 24	2.94 24 3.72 24 2.98 24	3.33 24 2.74 23 2.43 24	0.30 24 0.24 24 0.73 24	3.80 24	5.08 24	4.88 24	

Move the mouse over the numbers in the table to see the corresponding images. Click to compare with the ground truth.



(all, disc, and untext; see Sect. 4.3), resulting in a set of 24 scores per method. We sort each table by the average rank across all 24 scores to provide an ordering that *roughly* reflects the overall performance on the current metric and statistic. We want to emphasize that we do not aim to provide an overall ranking among the submitted methods. Authors sometimes report the rank of their method on one or more of the 32 sequences (often average angular error); however, many of the other 31 metric/statistic combinations might be better suited to compare the algorithms, depending on the application.

the ground-truth flows. Next to each score, the corresponding rank in the current column is indicated with a smaller blue number. The minimum (best) score in each column is shown in boldface. The methods are sorted by their average rank, which is computed over all 24 columns (eight sequences times three region masks each). The average rank serves as an *approximate* measure of performance *under the selected metric/statistic*.

cation of interest. Also note that the exact rank within any of the tables only gives a rough measure of performance, as there are various other ways that the scores across the 24 columns could be combined. We also list the runtimes reported by authors on the *Urban* sequence on the evaluation website (see Table 1). We made no attempt to normalize for the programming environment, CPU speed, number of cores, or other hardware acceleration. These numbers should be treated as a very rough guideline of the inherent computational complexity of the algorithms.

Table 1 Reported runtimes on the *Urban* sequence in seconds. We do not normalize for the programming environment, CPU speed, number of cores, or other hardware acceleration. These numbers should be treated as a very rough guideline of the inherent computational complexity of the algorithms

Algorithm	Runtime	Algorithm	Runtime
Adaptive (Wedel et al. 2009)	9.2	Seg OF (Xu et al. 2008)	60
Complementary OF (Zimmer et al. 2009)	44	Learning Flow (Sun et al. 2008)	825
Aniso. Huber-L1 (Werlberger et al. 2009)	2	Filter Flow (Seitz and Baker 2009)	34,000
DPOF (Lei and Yang 2009)	261	Graph Cuts (Cooke 2008)	1,200
TV-L1-improved (Wedel et al. 2008)	2.9	Black & Anandan (Black and Anandan 1996)	328
CBF (Trobis et al. 2008)	69	SPSA-learn (Li and Huttenlocher 2008)	200
Brox et al. (Brox et al. 2004)	18	Group Flow (Ren 2008)	600
Rannacher (Rannacher 2009)	0.12	2D-CLG (Bruhn et al. 2005)	844
F-TV-L1 (Wedel et al. 2008)	8	Horn & Schunck (Horn and Schunck 1981)	49
Second-order prior (Trobis et al. 2008)	14	TI-DOFE (Cassisa et al. 2009)	260
Fusion (Lempitsky et al. 2008)	2,666	FOLKI (Le Besnerais and Champagnat 2005)	1.4
Dynamic MRF (Glocke et al. 2008)	366	Pyramid LK (Lucas and Kanade 1981)	11.9

Table 2 A comparison of the average endpoint error (Avg. EE) results for 2D-CLG (Bruhn et al. 2005) (overall the best-performing algorithm in our preliminary study, Baker et al. 2007) and the best result uploaded to the evaluation website at the time of writing (Fig. 7)

	Army	Mequon	Schefflera	Wooden	Grove	Urban	Yosemite	Teddy
Best	0.09	0.18	0.24	0.18	0.74	0.39	0.08	0.50
2D-CLG (Bruhn et al. 2005)	0.28	0.67	1.12	1.07	1.23	1.54	0.10	1.38

Finally, we report on the evaluation website for each method the number of input frames and whether color information was utilized. At the time of writing, all of the 24 published methods discussed in this paper use only 2 frames as input; and 10 of them use color information.

The best-performing algorithm (both in terms of average endpoint error and average angular error) in our preliminary study (Baker et al. 2007) was 2D-CLG (Bruhn et al. 2005). In Table 2, we compare the results of 2D-CLG with the current best result in terms of average endpoint error (Avg. EE). The first thing to note is that performance has dramatically improved, with average EE values of less than 0.2 pixels on four of the datasets (*Yosemite*, *Army*, *Mequon*, and *Wooden*). The common elements of the more difficult sequences (*Grove*, *Teddy*, *Urban*, and *Schefflera*) are the presence of large motions and strong motion discontinuities. The complex discontinuities and fine structures of *Grove* seem to cause the most problems for current algorithms. A visual inspection of some computed flows (Fig. 8) shows that oversmoothing motion discontinuities is common even for the top-performing algorithms. A possible exception is DPOF (Lei and Yang 2009). On the other hand, the problems of complex non-rigid motion confounded with illumination changes, moving shadows, and real sensor noise (*Army*, *Mequon*, *Wooden*) do not appear to present as much of a problem for current algorithms.

5.2 Analysis of the Flow Errors

We now analyze the correlation between the metrics, statistics, region masks, and datasets for the flow errors. Figure 9 compares the average ranks computed over different subsets of the 32 pages of results, each of which contains 24 results for each algorithm. Column (a) contains the average rank computed over seven of the eight statistics (the standard deviation is omitted) and the three region masks for the endpoint error (EE). Column (b) contains the corresponding average rank for the angular error (AE). Columns (c) contain the average rank for each of the seven statistics for the endpoint error (EE) computed over the three masks and the eight datasets. Columns (d) contain the average endpoint error (Avg. EE) for each of the three masks just computed over the eight datasets. Columns (e) contains the Avg. EE computed for each of the datasets, averaged over each of the three masks. The order of the algorithms is the same as Fig. 7, i.e., we order by the average endpoint error (Avg. EE), the highlighted, leftmost column in (c). To help visualize the numbers, we color-code the average ranks with a color scheme where green denotes low values, yellow intermediate, and red large values.

We also include the Pearson product-moment coefficient r between various subsets of pairs of columns at the bottom of the figure. The Pearson measure of correlation takes

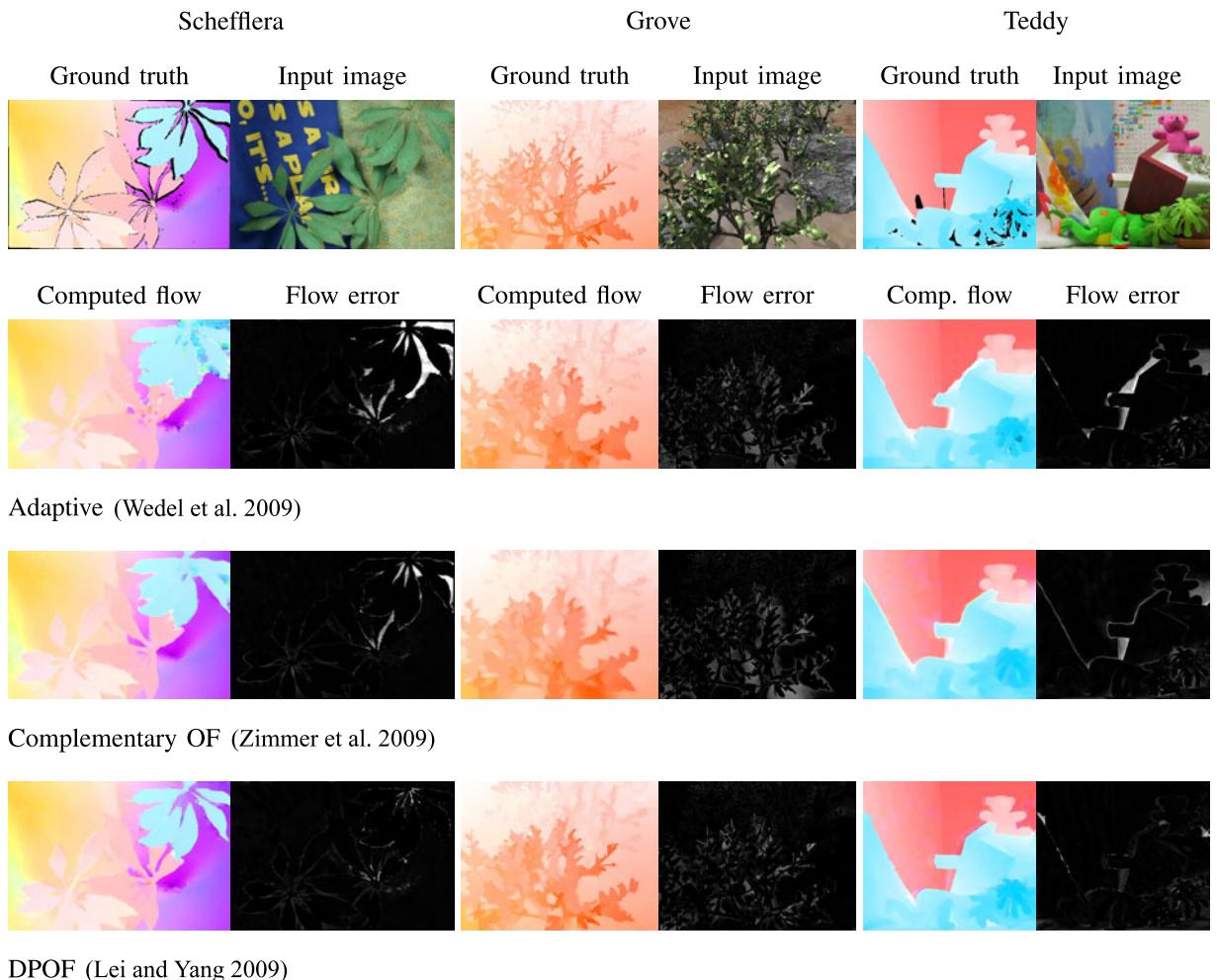


Fig. 8 The results of some of the top-performing methods on three of the more difficult sequences. All three sequences contain strong motion discontinuities. *Grove* also contains particularly fine structures.

The general tendency is to oversmooth motion discontinuities and fine structures. A possible exception is DPOF (Lei and Yang 2009)

on values between -1.0 and 1.0 , with 1.0 indicating perfect correlation. First, we include the correlation between each column and column (a). As expected, the correlation of column (a) with itself is 1.0 . We also include the correlation between all pairs of the statistics, between all pairs of the masks, and between all pairs of the datasets. The results are shown in the 7×7 , 3×3 , and 8×8 (symmetric) matrices at the bottom of the table. We color-code the correlation results with a separate scale where 1.0 is dark green and yellow/red denote lower values (less correlation).

5.2.1 Comparison of the Endpoint Error and the Angular Error

Columns (a) and (b) in Fig. 9 contain average ranks for the endpoint error (EE) and angular error (AE). The rankings generated with these two measures are highly correlated ($r = 0.989$), with only a few ordering reversals.

At first glance, it may seem that the two measures could be used largely interchangeably. Studying the qualitative results contained in Fig. 10 for the Complementary OF algorithm (Zimmer et al. 2009) on the Urban sequence leads to a different conclusion. The Complementary OF algorithm (which otherwise does very well) fails to correctly estimate the flow of the building in the bottom left. The average AE for this result is 4.64 degrees which ranks 6th in the table at the time of writing. The average EE is 1.78 pixels which ranks 20th at the time of writing. The huge discrepancy is due to the fact that the building in the bottom left has a very large motion, so the AE in that region is downweighted. Based on this example, we argue that the endpoint error (EE) should become the preferred measure of flow accuracy.

5.2.2 Comparison of the Statistics

Columns (c) in Fig. 9 contains a comparison of the various statistics, the average (Avg), the robustness mea-

Flow accuracy - analysis of statistics

Method	Avg over all stats		Individual EE statistics							Avg EE by mask			Avg EE by dataset								
	EE	AE	Avg	R0.5	R1.0	R2.0	A50	A75	A95	all	disc	untext	Army	Mequ.	Scheffl.	Wood.	Grove	Urban	Yosem.	Teddy	
Adaptive	4.8	4.4	4.4	5.2	5.3	4.7	3.9	3.7	6.7	4.3	4.9	4.0	1.0	4.7	7.0	1.7	3.7	3.0	10.7	3.3	
Complementary OF	5.2	5.9	5.7	4.3	3.7	3.8	6.5	6.0	6.4	6.1	3.9	7.1	5.7	1.0	2.3	3.3	9.0	13.7	5.3	5.3	
Aniso. Huber-L1	6.6	6.7	5.8	7.4	6.4	6.0	6.9	6.8	7.2	6.0	4.5	6.9	3.0	10.3	9.7	4.3	2.0	1.0	13.3	2.7	
DPOF	6.7	8.0	6.1	7.5	5.6	4.5	9.4	7.9	6.1	5.9	5.6	6.8	9.3	5.7	1.7	4.7	1.0	8.3	17.0	1.0	
TV-L1-improved	6.8	6.7	7.2	6.3	6.0	6.9	5.7	6.1	9.4	7.5	8.0	6.1	1.3	2.7	7.0	6.7	4.3	12.0	15.7	8.0	
CBF	8.1	9.2	7.8	8.6	7.3	6.8	7.9	8.6	9.4	8.0	6.6	8.6	3.3	10.3	6.0	6.7	5.0	1.7	21.0	8.0	
Brox et al.	8.7	8.4	8.4	9.6	9.8	7.8	7.6	8.2	9.4	7.6	8.8	8.8	8.3	9.3	5.0	8.7	14.0	7.7	2.3	11.7	
Rannacher	8.1	7.5	8.5	7.2	7.8	9.1	6.4	7.2	10.8	8.8	9.8	7.0	5.0	7.0	11.7	9.3	6.0	11.7	10.7	6.7	
F-TV-L1	8.3	8.1	8.8	6.8	7.9	10.5	6.8	7.5	9.5	8.8	9.4	8.4	13.3	11.7	10.7	13.3	5.0	4.7	8.0	4.0	
Second-order prior	9.6	10.2	9.0	11.0	11.4	9.5	6.9	8.8	10.5	8.9	9.8	8.3	5.0	8.3	11.3	4.3	8.7	7.7	16.3	10.0	
Fusion	9.3	11.8	9.4	9.1	8.7	8.2	9.9	8.8	11.0	8.3	9.9	10.1	8.0	2.0	3.3	7.7	12.3	9.3	17.3	15.3	
Dynamic MRF	10.3	9.7	11.1	9.2	9.9	10.8	8.8	9.9	12.7	10.4	12.3	10.8	11.0	5.0	5.0	10.3	14.3	16.7	8.3	18.3	
SegOF	12.2	12.4	11.7	14.8	13.7	9.5	12.8	13.1	10.0	12.8	10.3	12.0	12.3	16.3	13.7	10.7	17.0	16.0	2.3	5.0	
Learning Flow	12.6	11.7	13.3	11.7	13.0	14.0	9.5	12.6	14.0	12.6	15.8	11.6	5.7	11.0	11.0	15.7	18.7	16.0	13.3	15.3	
Filter Flow	14.2	14.2	14.3	14.5	14.4	13.2	15.2	15.5	12.6	14.8	13.9	14.3	15.3	14.0	17.0	18.0	15.3	4.7	18.7	11.3	
Graph Cuts	14.5	14.5	14.5	15.6	15.4	12.0	15.6	15.2	13.5	15.5	12.0	16.1	15.0	17.3	10.0	11.0	9.7	17.7	18.7	17.0	
Black & Anandan	15.5	15.7	15.0	15.2	15.8	16.7	15.2	15.6	14.8	15.1	16.0	13.8	17.3	16.7	17.7	17.3	12.7	12.7	10.7	14.7	
SPSA-learn	15.0	14.9	15.7	14.8	13.8	14.6	14.8	14.5	16.5	15.8	15.1	16.1	17.3	16.7	17.0	18.3	14.3	18.0	5.7	18.0	
Group Flow	16.2	16.3	15.9	17.4	18.3	15.5	16.2	16.6	13.5	16.5	15.8	15.5	19.0	21.0	18.3	12.0	17.0	20.3	6.0	13.7	
2D-CLG	17.3	15.9	17.4	18.3	18.2	17.5	16.8	17.2	15.5	17.4	16.6	18.1	20.7	18.7	21.3	21.7	19.0	15.7	2.3	19.7	
Horn & Schunck	18.6	19.1	18.6	18.8	19.1	19.0	18.9	19.0	16.5	18.1	20.0	17.6	20.3	19.0	20.0	20.0	19.7	17.0	11.7	21.0	
TI-DOFE	19.8	20.7	19.6	21.2	20.8	19.2	20.9	20.0	17.2	18.6	20.5	19.6	23.0	22.0	23.3	23.0	20.7	16.3	6.3	22.0	
FOLKI	22.2	21.8	22.6	22.1	22.5	21.8	21.5	22.3	22.6	22.4	23.1	22.4	22.7	23.3	22.0	21.3	23.0	23.0	22.7	23.0	
Pyramid LK	23.2	23.1	23.7	22.8	23.3	23.4	23.2	23.3	22.9	24.0	23.1	24.0	23.0	23.7	23.3	24.0	24.0	23.7	24.0	24.0	
Correlation with EE:	1.0	.989	.996	.985	.989	.977	.973	.993	.954	.992	.971	.986	.919	.913	.899	.920	.879	.755	.158	.870	
Correlation in group:			Avg	R0.5	R1.0	R2.0	A50	A75	A95	all	disc	untext	Army	Mequ.	Scheffl.	Wood.	Grove	Urban	Yosem.	Teddy	
			Avg	1.0	.970	.978	.981	.962	.986	.968	.960	.985	.910	.873	.827	.887	.783	.693	-.045	.759	
			R0.5	.970	1.0	.991	.937	.972	.985	.903	.960	1.0	.937	.873	1.0	.909	.831	.725	.595	.041	.658
			R1.0	.978	.991	1.0	.962	.952	.980	.922	.985	.937	1.0	.827	.909	1.0	.878	.741	.569	.005	.675
			R2.0	.981	.937	.962	1.0	.915	.954	.964	.962	.972	.951	.887	.831	.831	.634	.026	.424	.813	
			A50	.962	.972	.952	.915	1.0	.986	.888	.962	.972	.951	.783	.725	.741	.766	1.0	-.042	.744	
			A75	.986	.985	.980	.954	.986	1.0	.925	.968	.903	.922	.964	.888	.925	1.0	.766	.041	.135	
			A95	.968	.903	.922	.964	.888	.925	1.0	.968	.903	.922	.964	.888	.925	.783	.852	.744	.135	

(a) (b)

(c)

(d)

(e)

Fig. 9 A comparison of the various different metrics, statistics, region masks, and datasets for flow errors. Each column contains the average rank computed over a different subset of the 32 pages of results, each of which contains 24 different results for each algorithm. See the main body of the text for a description of exactly how each column is computed. To help visualize the numbers, we color-code the average ranks with a color scheme where green denotes low values, yellow intermediate, and red large values. The order of the algorithms is the

same as Fig. 7, i.e., we order by the average endpoint error (Avg. EE), the leftmost column in (e), which is highlighted in the table. At the bottom of the table, we include correlations between various subsets of pairs of the columns. Specifically, we compute the Pearson product-moment coefficient r . We separately color-code the correlations with a scale where dark green is 1.0 and yellow/red denote lower values

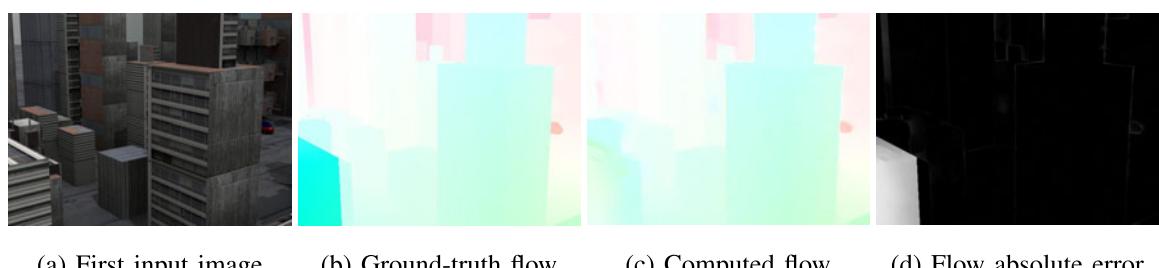


Fig. 10 Results of the Complementary OF algorithm (Zimmer et al. 2009) on the Urban sequence. The average AE is 4.64 degrees which ranks 6th in the table at the time of writing. The average EE is 1.78 pixels which ranks 20th at the time of writing. The huge discrepancy is

due to the fact that the building in the bottom left has a very large motion, so the AE in that region is downweighted. Based on this example, we argue that the endpoint error (EE) should become preferred measure of flow accuracy

sures (R0.5, R1.0, and R2.0), and the accuracy measures (A50, A75, and A95). The first thing to note is that again these measures are all highly correlated with the average over all the statistics in column (a) and with themselves.

The outliers and variation in the measures for any one algorithm can be very informative. For example, the performance of DPOF (Lei and Yang 2009) improves dramatically from R0.5 to R2.0 and similarly from A50 to A95.

This trend indicates that DPOF is good at avoiding gross outliers but is relatively weak at obtaining high accuracy. DPOF (Lei and Yang 2009) is a segmentation-based discrete optimization algorithm, followed by a continuous refinement (Sect. 2.4.2). The variation of the results across the measures indicates that the combination of segmentation and discrete optimization is beneficial in terms of avoiding outliers, but that perhaps the continuous refinement is not as sophisticated as recent purely continuous algorithms. The qualitative results obtained by DPOF on the Schefflera and Grove sequences in Fig. 8 show relatively good results around motion boundaries, supporting this conclusion.

5.2.3 Comparison of the Region Masks

Columns (d) in Fig. 9 contain a comparison of the region masks, *All*, *Disc*, and *Untext*. Overall, the measures are highly correlated by rank, particularly for the *All* and *Untext* masks. When comparing the actual error scores in the individual tables (e.g., Fig. 7), however, the errors are much higher throughout in the *Disc* regions than in the *All* regions, while the errors in the *Untext* regions are typically the lowest. As expected, the *Disc* regions thus capture what is still the hardest task for optical flow algorithms: to accurately recover motion boundaries. Methods that strongly smooth across motion discontinuities (such as the Horn and Schunck algorithm 1981, which uses a simple L2 prior) also show a worse performance for *Disc* in the rankings (columns (d) in Fig. 9). Textureless regions, on the other hand, seem to be no problem for today's methods, essentially all of which optimize a global energy.

5.2.4 Comparison of the Datasets

Columns (e) in Fig. 9 contain a comparison across the datasets. The first thing to note is that the results are less strongly correlated than across statistics or region masks. The results on the *Yosemite* sequence, in particular, are either poorly or negatively correlated with all of the others. (The main reason is that the *Yosemite* flow contains few discontinuities and consequently methods do well here that over-smooth other sequences with more motion boundaries.) The most correlated subset of results appear to be the four hidden texture sequences *Army*, *Mequon*, *Schefflera*, and *Wooden*. These results show how performance on any one sequence can be a poor predictor of performance on other sequences and how a good benchmark needs to contain as diverse a set of data as possible. Conversely, any algorithm that performs consistently well across a diverse collection of datasets can probably be expected to perform well on most inputs.

Studying the results in detail, a number of interesting conclusions can be noted. Complementary OF (Zimmer et al. 2009) does well on the hidden texture data (*Army*,

Mequon, *Schefflera*, *Wooden*) presumably due to the use of a relatively sophisticated data term, including the use of a different robust penalization function for each channel in HSV color space (the hidden texture data contains a number of moving shadows and other illumination-related effects), but not as well on the sequences with large motion (*Urban*) and complex discontinuities (*Grove*). DPOF (Lei and Yang 2009), which involves segmentation and performs best on *Grove*, does particularly poorly on *Yosemite* presumably because segmenting the grayscale *Yosemite* sequence is difficult. F-TV-L1 (Wedel et al. 2008) does well on the largely rigid sequences (*Grove*, *Urban*, *Yosemite*, and *Teddy*), but poorly on the non-rigid sequences (*Army*, *Mequon*, *Schefflera*, and *Wooden*). F-TV-L1 uses a rigidity prior and so it seems that this component is being used too aggressively. Note, however, that a later algorithm by the same group of researchers (Adaptive—Wedel et al. 2009—which also uses a rigidity prior) appears to have addressed this problem. The flow fields for Dynamic MRF (Glocker et al. 2008) all appear to be over-smoothed; however, quantitatively, the performance degradation is only apparent on the sequences with strong discontinuities (*Grove*, *Urban*, and *Teddy*). In summary, the relative performance of an algorithm across the various datatypes in our benchmark can lead to insights into which of its components work well and which are limiting performance.

5.3 Analysis of the Interpolation Errors

We now analyze the correlation between the metrics, statistics, region masks, and datasets for the interpolation errors. In Fig. 11, we include results for the interpolation errors that are analogous to the flow error results in Fig. 9, described in Sect. 5.2. Note that we are now comparing interpolated frames (generated from the submitted flow fields using the interpolation algorithm from Sect. 3.3.2) with the true intermediate frames. Also, recall that we use a different set of test sequences for the interpolation evaluation: the four high-speed datasets *Backyard*, *Basketball*, *Dumptruck*, and *Evergreen*, in addition to *Mequon*, *Schefflera*, *Urban*, and *Teddy*, as representatives of the three other types of datasets. We sort the algorithms by the average interpolation error performance (Avg. IE), the leftmost column in Fig. 11(c). The ordering of the algorithms in Fig. 11 is therefore different from that in Fig. 9.

5.3.1 Comparison of the Interpolation and Normalized Interpolation Errors

Columns (a) and (b) in Fig. 11 contain average ranks for the interpolation error (IE) and the normalized interpolation error (NE). The rankings generated with these two measures are highly correlated ($r = 0.981$), with only a few ordering

Interpolation accuracy - analysis of statistics

Method	Avg over all stats		Individual IE statistics							Avg IE by mask			Avg IE by dataset								
	IE	NE	Avg	R2.5	R5.0	R10	A90	A95	A99	all	disc	untext	Mequ.	Scheffl.	Urban	Teddy	Backyd.	Basktb.	Dumpr.	Evergr.	
CBF	5.6	7.9	3.5	10.6	7.4	3.8	4.6	5.0	4.0	2.4	2.3	6.0	2.3	5.3	1.3	3.3	4.0	3.0	3.7	5.3	
Aniso. Huber-L1	4.5	4.0	4.6	5.8	5.5	4.2	3.2	3.8	4.1	4.1	4.6	5.0	4.0	11.3	2.3	4.0	8.3	1.7	1.0	4.0	
Second-order prior	5.2	4.4	5.5	4.9	5.1	6.1	4.0	5.0	5.8	5.1	6.4	4.9	3.3	8.0	6.0	3.0	6.3	4.3	3.0	9.7	
Brox et al.	4.6	4.8	6.3	5.4	4.5	3.8	3.0	3.8	5.2	6.3	6.4	6.3	5.7	3.0	4.7	3.3	2.3	14.0	16.3	1.0	
F-TV-L1	6.5	8.0	7.1	9.0	7.4	5.4	4.0	5.8	7.1	5.8	6.1	9.4	14.7	11.0	5.0	7.7	4.0	2.7	5.7	6.0	
Filter Flow	11.7	12.8	9.7	14.8	13.1	10.8	10.5	12.0	10.7	8.9	8.6	11.5	10.7	16.0	9.0	9.3	5.3	9.7	7.0	10.3	
Fusion	9.3	8.3	10.0	7.3	10.5	11.6	8.2	8.1	9.2	10.3	10.1	9.8	4.7	2.0	6.3	6.7	13.3	21.3	10.0	16.0	
Black & Anandan	11.0	10.4	10.1	12.9	12.7	11.1	10.0	9.9	10.2	11.1	9.5	9.6	12.7	17.7	15.7	12.3	4.0	7.7	7.7	3.0	
DPOF	9.7	10.1	10.2	11.8	10.4	8.8	8.1	9.2	10.0	9.9	12.0	8.6	15.0	1.0	15.3	6.7	13.7	9.0	8.7	12.0	
2D-CLG	10.4	10.9	11.0	11.5	11.7	11.0	8.7	8.3	10.6	10.5	8.1	14.4	8.0	15.7	9.7	12.3	6.7	13.3	5.0		
Horn & Schunck	13.9	13.9	11.1	17.3	16.4	14.0	13.1	13.0	12.2	12.1	10.8	10.5	9.0	20.0	13.7	16.3	4.7	5.3	13.0	7.0	
Adaptive	10.1	9.3	12.5	9.1	10.2	10.2	7.7	9.6	11.4	13.1	14.9	9.5	11.7	16.7	7.0	12.0	14.3	14.3	12.0	12.0	
Complementary OF	9.8	8.8	12.5	6.1	7.6	11.6	7.8	9.8	13.4	14.0	14.5	8.9	13.3	4.3	19.0	13.0	14.7	9.0	11.0	15.3	
TV-L1-improved	11.0	11.6	12.8	10.8	11.3	10.9	7.8	10.7	12.5	13.1	12.9	12.4	8.3	15.3	11.0	5.0	11.7	18.3	14.3	14.3	
Graph Cuts	10.2	11.5	13.0	7.1	8.0	10.9	9.0	11.0	12.2	14.1	12.9	11.9	17.0	5.3	14.0	12.0	15.7	10.3	15.7	13.7	
TI-DOFE	14.7	14.4	13.5	17.4	16.2	14.4	13.8	13.9	13.4	13.4	12.9	14.4	17.3	22.3	9.0	19.0	5.0	12.7	7.3	15.7	
Dynamic MRF	14.5	14.5	14.5	12.8	14.5	16.2	12.4	14.8	16.2	14.0	15.0	14.6	9.0	7.3	12.7	18.0	12.3	22.0	18.3	16.7	
Learning Flow	17.3	18.1	15.8	20.1	20.0	17.6	15.8	16.0	15.7	15.6	15.6	16.3	11.0	13.3	24.0	13.7	19.3	15.3	12.0	18.0	
FOLKI	19.0	20.9	15.9	22.5	22.2	18.6	19.0	18.9	16.2	14.3	13.9	19.5	20.3	22.7	15.0	20.7	11.3	16.3	10.7	10.0	
Rannacher	12.6	13.2	16.0	11.5	11.9	12.5	9.0	12.0	15.0	16.9	17.5	13.8	13.0	18.0	15.3	13.0	15.7	19.0	18.3	16.0	
SPSA-learn	14.5	15.6	18.0	10.7	12.6	15.2	12.4	14.9	17.8	18.8	18.4	16.9	19.0	12.3	21.7	18.3	17.7	12.3	24.0	18.7	
SegOf	14.3	14.0	18.1	11.7	12.0	15.2	11.9	13.8	17.5	19.1	18.8	16.5	19.0	7.7	19.3	22.0	21.3	19.3	21.7	14.7	
Group Flow	19.8	18.9	21.1	18.3	19.8	21.7	19.0	18.4	20.2	21.1	21.8	20.4	23.0	15.7	20.3	22.3	23.0	23.7	18.7	22.0	
Pyramid LK	21.9	22.3	22.2	21.9	21.9	22.5	22.5	20.9	21.8	21.6	22.9	21.3	22.4	23.3	23.7	22.7	24.0	22.0	15.7	22.0	
Correlation with IE:	1.0	.981	.916	.876	.956	.983	.987	.991	.937	.874	.835	.946	.766	.610	.796	.901	.621	.605	.592	.725	
Correlation in group:			Avg	R2.5	R5.0	R10	A90	A95	A99	all	disc	untext	Mequ.	Scheffl.	Urban	Teddy	Backyd.	Basktb.	Dumpr.	Evergr.	
			1.0	633	769	933	859	918	988	1.0	.979	.877	.408	.422	.725	.819	.564	.395	.539	.599	
			R2.5	.633	1.0	.963	.793	.900	.847	.670	.979	1.0	.824	.422	.262	1.0	.549	.065	.091	.142	.173
			R5.0	.769	.963	1.0	.914	.968	.929	.799	.877	.824	1.0	.725	.262	1.0	.738	.238	.459	.658	.687
			R10	.933	.793	.914	1.0	.968	.973	.954	.918	.976	.885	.564	.065	.723	.819	.549	.611	.622	
			A90	.859	.900	.968	.968	1.0	.976	.885	.918	.946	.946	.395	.091	.459	.507	.582	1.0	.692	.680
			A95	.918	.847	.929	.973	.976	1.0	.946	.988	.670	.799	.687	.622	.759	.680	.596	1.0	.596	
			A99	.988	.670	.799	.954	.885	.946	1.0	.988	.946	.946	.599	.173	.687	.622	.759	.680	.596	1.0

(a) (b)

(c)

(d)

(e)

Fig. 11 A comparison of the various different metrics, statistics, region masks, and datasets for interpolation errors. These results are analogous to those in Fig. 9, except the results here are for interpolation errors rather than flow errors. See Sect. 5.2 for a description of

how this table was generated. We sort the algorithms by the average interpolation error performance (Avg. IE), the first column in (c). The ordering of the algorithms is therefore different to that in Fig. 9

reversals. Most of the differences between the two measures can be explained by the relative weight given to the discontinuity and textureless regions. The rankings in columns (a) and (b) are computed by averaging the ranking over the three masks. The normalized interpolation error (NE) generally gives additional weight to textureless regions, and less weight to discontinuity regions (which often also exhibit an intensity gradient). For example, CBF (Trobin et al. 2008) performs better on the *All* and *Disc* regions than it does on the *Untext* regions, which explains why the NE rank for this algorithm is slightly higher than the IE rank.

5.3.2 Comparison of the Statistics

Columns (c) in Fig. 11 contain a comparison of the various statistics, the average (Avg), the robustness measures (R2.5, R5.0, and R10.0), and the accuracy measures (A90, A95, and A99). Overall the results are highly correlated. The most obvious exception is R2.5, which measures the percentage of pixels that are predicted very precisely (within 2.5 graylevels). In regions with some texture, very accurate flow is needed to obtain the highest possible precision.

Algorithms such as CBF (Trobin et al. 2008) and DPOF (Lei and Yang 2009), which are relatively robust but not so accurate (compare the performance of these algorithms for R0.5 and R10.0 in Fig. 9), therefore perform worse in terms of R2.5 than they do in terms of R5.0 and R10.0.

5.3.3 Comparison of the Region Masks

Columns (d) in Fig. 11 contain a comparison of the region masks, *All*, *Disc*, and *Untext*. The *All* and *Disc* results are highly correlated, whereas the *Untext* results are less correlated with the other two masks. Studying the detailed results on the webpage for the outliers in columns (d), there does not appear to be any obvious trend. The rankings in the *Untext* regions just appear to be somewhat more “noisy” due to the fact that for some datasets there are relatively few *Untext* pixels and all algorithms have relatively low interpolation errors in those regions. The actual error values (as opposed to their rankings) are quite different between the three regions masks. Like the flow accuracy errors (Sect. 5.2.3), the IE values are highest in the *Disc* regions since flow errors near object boundaries usually cause interpolation errors as well.

5.3.4 Comparison of the Datasets

Columns (e) in Fig. 11 contain a comparison across the datasets. The results are relatively uncorrelated, just like the flow errors in Fig. 9. The most notable outlier for interpolation is *Schefflera*. Studying the results in detail on the website, the primary cause appears to the right hand side of the images, where the plant leaves move over the textured cloth. This region is difficult for many flow algorithms because the difference in motions is small and the color difference is not great either. Only a few algorithms (e.g., DPOF—Lei and Yang 2009, Fusion—Lempitsky et al. 2008, and Dynamic MRF—Glocker et al. 2008) perform well in this region. Getting this region correct is more important in the interpolation study than in the flow error study because: (1) the background is quite highly textured, so a small flow error leads to a large interpolation error (see the error maps on the webpage) and (2) the difference between the foreground and background flows is small, so oversmoothing the foreground flow is not penalized by a huge amount in the flow errors. The algorithms that perform well in this region do not perform particularly well on the other sequences, as none of the other seven interpolation datasets contain regions with similar causes of difficulty, leading to the results being fairly uncorrelated.

5.4 Comparison of the Flow and Interpolation Errors

In Fig. 12, we compare the flow errors with the interpolation errors. In the left half of the figure, we include the average rank scores, computed over all statistics (except the standard deviation) and all three masks. We compare flow endpoint errors (EE), interpolation errors (IE), and normalized interpolation errors (NE), and include two columns for each, Avg and Avg4. The first column, Avg EE, is computed over all eight flow error datasets, and corresponds exactly to column (a) in Fig. 9. Similarly, the third and fifth columns, Avg IE and Avg NE, are computed over all eight interpolation error datasets, and correspond exactly to columns (a) and (b) in Fig. 11. To remove any dependency on the different datasets, we provide the Avg4 columns, which are computed over the four sequences that are common to the flow and interpolation studies: *Mequon*, *Schefflera*, *Urban*, and *Teddy*.

The right half of Fig. 12 shows the 6×6 matrix of the column correlations. It can be seen that the correlation between the results for Avg4 EE and Avg4 IE is only 0.763. The comparison here uses the same datasets, statistics, and masks; the only difference is the error metric, flow endpoint error (EE) vs. interpolation error (IE). Part of the reason these measures are relatively uncorrelated is that the

Correlation of flow and interpolation accuracy

Method	EE		IE		NE		Correlation:	EE		IE		NE	
	Avg	Avg4	Avg	Avg4	Avg	Avg4		Avg	Avg4	Avg	Avg4	Avg	Avg4
Adaptive	4.4	4.5	12.5	11.8	9.8	10.4							
Complementary OF	5.7	5.6	12.5	12.4	11.0	9.3							
Aniso. Huber-L1	5.8	5.9	4.6	5.4	5.0	5.1							
DPOF	6.1	4.2	10.2	9.5	10.9	10.3							
TV-L1-improved	7.2	7.4	12.8	9.9	12.7	9.8							
CBF	7.8	6.5	3.5	3.1	5.6	4.8							
Brox et al.	8.4	8.4	6.3	4.2	7.5	4.8							
Rannacher	8.5	9.3	16.0	14.8	14.1	13.2							
F-TV-L1	8.8	7.8	7.1	9.6	8.4	9.2							
Second-order prior	9.0	9.3	5.5	5.1	5.5	5.1							
Fusion	9.4	7.5	10.0	4.9	8.7	6.3							
Dynamic MRF	11.1	11.3	14.5	11.8	15.3	11.3							
SegOF	11.7	12.8	18.1	17.0	15.3	15.8							
Learning Flow	13.3	13.3	15.8	15.5	15.2	15.6							
Filter Flow	14.3	11.8	9.7	11.3	11.0	14.0							
Graph Cuts	14.5	15.5	13.0	12.1	13.0	11.8							
Black & Anandan	15.0	15.4	10.1	14.6	10.1	14.5							
SPSA-learn	15.7	17.4	18.0	17.8	19.0	18.4							
Group Flow	15.9	18.3	21.1	20.3	19.2	18.8							
2D-CLG	17.4	18.8	11.0	11.4	11.6	11.3							
Horn & Schunck	18.6	19.3	11.1	14.8	10.4	14.0							
TI-DOFE	19.6	20.9	13.5	16.9	12.0	16.1							
FOLKI	22.6	22.8	15.9	19.7	18.0	19.8							
Pyramid LK	23.7	23.7	22.2	23.4	21.5	23.1							

Fig. 12 A comparison of the flow errors, the interpolation errors, and the normalized interpolation errors. We include two columns for the average endpoint error. The leftmost (Avg EE) is computed over all eight flow error datasets. The other column (Avg4 EE) is computed over the four sequences that are common to the flow and interpolation studies (*Mequon*, *Schefflera*, *Urban*, and *Teddy*). We also include two columns each for the average interpolation error and the average normalized interpolation error. The leftmost of each pair (Avg IE and

Avg NE) are computed over all eight interpolation datasets. The other columns (Avg4 IE and Avg NE) are computed over the four sequences that are common to the flow and interpolation studies (*Mequon*, *Schefflera*, *Urban*, and *Teddy*). On the right, we include the 6×6 matrix of the correlations of the six columns on the left. As in previous figures, we separately color-code the average rank columns and the 6×6 correlation matrix

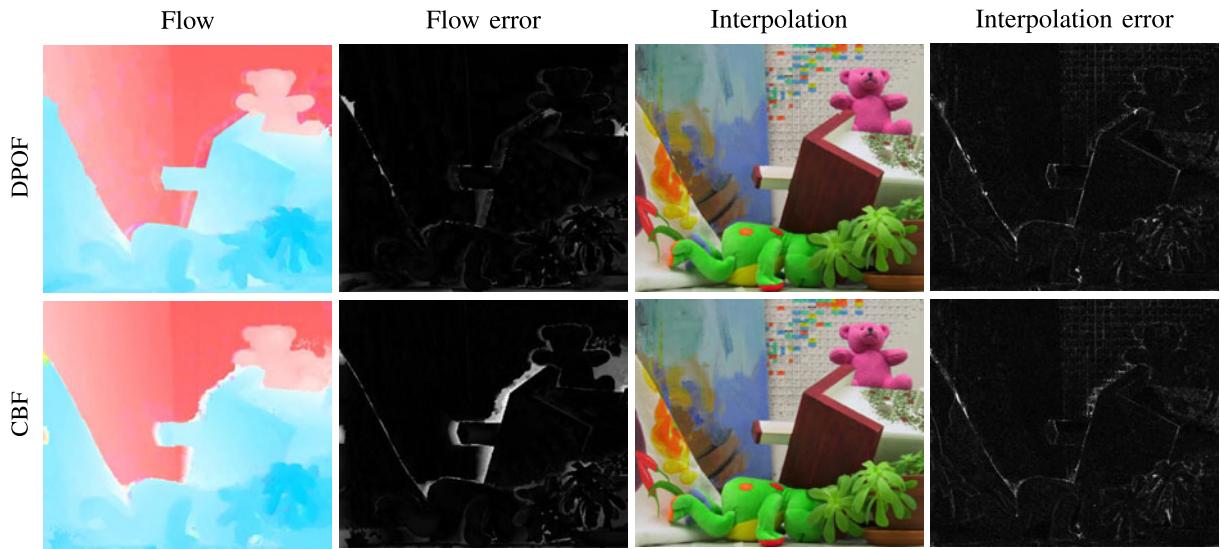


Fig. 13 A comparison of the flow and interpolation results for DPOF (Lei and Yang 2009) and CBF (Trobin et al. 2008) on the *Teddy* sequence to illustrate the differences between the two measures of performance. DPOF obtains the best flow results with an Avg. EE of

0.5 pixels, whereas CBF is ranked 9th with an Avg. EE of 0.76 pixels. CBF obtains the best interpolation error results with an Avg. IE of 5.21 graylevels, whereas DPOF is ranked 6th with an Avg. IE of 5.58 graylevels

interpolation errors are themselves a little noisy internally. As discussed above, the R2.5 and *Untext* mask results are relatively uncorrelated with the results for the other measures and masks. The main reason, however, is that the interpolation penalizes small flow errors in textured regions a lot, and larger flow errors in untextured regions far less. An illustration of this point is included in Fig. 13. We include both flow and interpolation results for DPOF (Lei and Yang 2009) and CBF (Trobin et al. 2008) on the *Teddy* sequence. DPOF obtains the best flow results with an average endpoint error of 0.5 pixels, whereas CBF is the 9th best with an average endpoint error of 0.76 pixels. CBF obtains the best interpolation error results with an average interpolation error of 5.21 graylevels, whereas DPOF is 6th best with an average interpolation error of 5.58 graylevels. Although the flow errors for CBF are significantly worse, the main errors occur where the foreground flow is “fattened” into the relatively textureless background to the left of the birdhouse and the right of the teddy bear. The interpolation errors in these regions are low. On the other hand, DPOF makes flow errors on the boundary between the white cloth and blue painting that leads to large interpolation errors. The normalized interpolation error (NE) is meant to compensate for this difference between the flow and interpolation errors. Figure 12 does show that the Avg4 NE and Avg4 EE measures are more correlated ($r = 0.803$) than the Avg4 IE and Avg4 EE measures ($r = 0.763$). The increased degree of correlation is marginal, however, due to the difficulty in setting a spatial smoothing radius for the gradient computation, and the need to regularize the NE measure by adding ϵ to the denominator. Therefore, as one might

expect, the performance of a method in the interpolation evaluation yields only limited information about the accuracy of the method in terms of recovering the true motion field.

5.5 Analysis of the Algorithms

Table 3 contains a summary of most of the algorithms for which results have been uploaded to our online evaluation. We omit the unpublished algorithms and a small number of the algorithms that are harder to characterize in terms of our taxonomy. We list the algorithms in the same order as Figs. 7 and 9. Generally speaking, the better algorithms are at the top, although note that this is just one way to rank the algorithms. For each algorithm, we mark which elements of our taxonomy in Sect. 2 it uses. In terms of the data term, we mark whether the algorithm uses the L1 norm or a different robust penalty function (Sect. 2.1.2). Neither column is checked for an algorithm such as Horn and Schunck (1981), which uses the L2 norm. We note if the algorithm uses a gradient component in the data term or any other more sophisticated features (Sect. 2.1.3). We also note if the algorithm uses an explicit illumination model (Sect. 2.1.4), normalizes the data term in any way, or uses a sophisticated color model to reduce the effects of illumination variation (Sect. 2.1.5).

For the spatial prior term, we also mark whether the algorithm uses the Total Variation (TV) norm or a different robust penalty function (Sect. 2.2.2). We note if the algorithm spatially weights the prior (Sect. 2.2.3) or if the weighting is anisotropic (Sect. 2.2.4). We also note if the algorithm

Table 3 A classification of most of the algorithms for which results have been uploaded to our online evaluation in terms of which elements of our taxonomy in Sect. 2 they use

Algorithm	Data term				Prior term				Optimization			Misc.					
	L1 norm	Other robust penalty F_n	Gradient/other features	Illum. modeling/norm/color	L1/TV norm	Other robust penalty F_n	Spatial weighting	Anisotropic weighting	Higher-order prior	Rigidity prior	Cont.-gradient descent	Cont.-variational/extremal	Cont.-other	Discr.-fusion	Discr.-reparameterization	Learning	Visibility/occlusion
Adaptive (Wedel et al. 2009)	X		X	X	X	X	X	X	X			X					
Complementary OF (Zimmer et al. 2009)	X	X	X			X	X	X				X					X
Aniso. Huber-L1 (Werlberger et al. 2009)	X		X			X	X	X					X				
DPOF (Lei and Yang 2009)			X	X		X					X			X		X	X
TV-L1-improved (Wedel et al. 2008)	X		X	X									X				
CBF (Trobin et al. 2008)	X				X				X				X				X
Brox et al. (Brox et al. 2004)	X	X		X								X					X
F-TV-L1 (Wedel et al. 2008)	X				X				X				X				
Second-order prior (Trobin et al. 2008)	X								X				X				
Fusion (Lempitsky et al. 2008)		X	X			X	X				X		X				X
Dynamic MRF (Glockner et al. 2008)	X				X									X			
Seg OF (Xu et al. 2008)	X				X	X						X				X	X
Learning Flow (Sun et al. 2008)		X	X			X	X				X				X		
Filter Flow (Seitz and Baker 2009)	X		X	X		X		X	X				X				X
Graph Cuts (Cooke 2008)	X				X									X			X
Black & Anandan (Black and Anandan 1996)		X				X					X						
SPSA-learn (Li and Huttenlocher 2008)		X				X					X				X		X
Horn & Schunck (Horn and Schunck 1981)											X						

uses a higher-order prior (Sect. 2.2.5) or a rigidity prior (Sect. 2.2.6).

In terms of the optimization algorithm, we mark if the algorithm uses a gradient-descent based continuous optimization (Sect. 2.3.1). We also specify which algorithms are variational or use other extremal approaches (Sect. 2.3.2). Other approaches (Sect. 2.3.3), such as the dual variable approach and the use of Linear Programming, are grouped together. In terms of discrete optimization, we distinguish fusion based algorithms (Sect. 2.4.1) from reparameterization based algorithms (Sect. 2.4.1) and note which approaches also use a continuous optimization phase to refine the results (Sect. 2.4.3).

Finally, we also denote which algorithms use learning (Sect. 2.5.1) to optimize the parameters and which algorithms perform explicit visibility or occlusion reasoning (Sect. 2.5.5). In the last column we mark whether the algorithm uses color images.

Based on Table 3, we note the following:

- *Degree of Sophistication:* The algorithms toward the top of the table tend to use a lot more of the refinements to the data and prior terms. Spatial weighting, anisotropic weighting, and the addition of robustness to illumination changes through data term normalization or the use of features, are all common components in the top-performing algorithms.

- *Choice of Penalty Function:* The L1 norm is a very popular choice, particularly for the data term. A couple of the top-performing algorithms combine a L1 norm on the data term with a different (more truncated) robust penalty function on the prior term.
- *Rigidity:* As discussed in Sect. 5.2.4, one algorithm that uses rigidity (F-TV-L1—Wedel et al. 2008) does poorly on the non-rigid scenes, however, Adaptive (Wedel et al. 2009) (a subsequent algorithm by the same researchers) does well on all sequences.
- *Continuous Optimization:* The gradient descent algorithms (discounting the ones that first perform a discrete optimization) all appear at the bottom of the table. On the other hand, the variational approaches appear throughout the table. Note that there is a correlation between the use of variational methods and more sophisticated energy functions that is not intrinsic to the variational approach. A direct comparison of different optimization methods with the same objective functions needs to be carried out. The dual-variable approach is competitive with the best algorithms, and may offer a speed advantage.
- *Discrete Optimization:* The discrete optimization algorithms do not perform particularly well. Note, however, that the energy functions used in these methods are generally relatively simple and might be extended in the future to incorporate some of the more sophisticated elements. It does, however, appear that refining the results with a continuous optimization is required to obtain good results (if accuracy is measured using average endpoint error).
- *Miscellaneous:* There are few algorithms that employ learning in the table, making it difficult to draw conclusions in terms of performance. This is likely to change in the future, as learning techniques are maturing and more labeled training data is becoming available. Similarly, few algorithms incorporate explicit visibility or occlusion reasoning, making it difficult to assess how important this could be. Notably, all 24 algorithms considered here utilize only 2 input frames, despite the fact that we make 8-frame sequences available. In contrast, on previous evaluation sets (particularly *Yosemite*) multi-frame methods relying on temporal smoothing were quite common. This raises the question of whether temporal smoothing, at least as applied so far, is less suited for the more challenging sequences considered here. A definitive answer to this point cannot be given in this paper, but should be subject of future work. Finally, less than half of the algorithms utilize color information, and there is no obvious correlation with performance. The utility of color for image matching clearly deserves further study as well; see Bleyer and Chambon (2010) for some recent insights on this issue in the context of stereo matching.

5.6 Comparison with State-of-the-Art Stereo Methods

As mentioned in Sect. 3.4, evaluating the flow algorithms on the modified Teddy stereo dataset allows a comparison with current stereo methods from the online Middlebury stereo evaluation at <http://vision.middlebury.edu/stereo/> (Scharstein and Szeliski 2002). To compare the state of the art, we select the best-performing flow and stereo methods from the two evaluations and compute the median of the lowest five R0.5 and R1.0 endpoint error scores on the Teddy dataset. Recall that the RX endpoint error score measures the percentage of pixels whose endpoint (or disparity) error is greater than X pixels. We compute these scores for both *All* and *Disc* region masks. While there are slight differences in the definition of the *Disc* region masks between flow and stereo evaluations, the comparison provides a good sense of the relative accuracy of the two classes of methods.

When comparing these scores, it becomes clear that the current top stereo methods significantly outperform the top flow methods. In particular, the median of the lowest five R1.0 error rates in *All* regions is 9.9 for flow, but only 6.5 for stereo (a reduction by 34%). In the *Disc* regions, the errors are much higher and the difference is even more pronounced, with a median error of 27.6 for flow and 10.0 for stereo (a reduction by 64%). Of course, stereo methods solve an easier problem, since correspondences are restricted to lie on epipolar lines, which may be one reason for the performance difference (though, as mentioned earlier, some flow methods employ rigidity priors that aid in the recovery of static scenes). Another significant difference is that many current stereo methods employ either discrete label sets to model disparities, or piecewise planar surface models. In contrast, current flow methods typically perform a continuous optimization. This explains why current stereo methods are able to recover much sharper depth discontinuities than most current flow methods, which is apparent both quantitatively from the *Disc* scores and qualitatively from examining the recovered disparity maps and flow fields.

When comparing the R0.5 scores, which reflect subpixel accuracy, the errors are higher overall, but the difference between the top stereo and flow methods is slightly less pronounced: the median of the lowest five scores in *All* regions is now 16.6 for flow, and 13.8 for stereo (a reduction by 17%); in the *Disc* regions the median is now 38.0 for flow, and 22.5 for stereo (a reduction by 41%). A possible explanation for the smaller performance difference when using the R0.5 scores is that the continuous approaches used in optical flow techniques are better able to achieve subpixel precision.

In summary, current flow algorithms, when run on a stereo pair, cannot quite match the performance of state-of-the-art stereo methods, particularly near depth discontinuities. Conversely, most current stereo methods use discrete

label sets or simplified surface models that cannot be easily adapted to the problem of recovering continuous and smoothly varying 2D motion fields. It is likely that stereo and flow algorithms will become more similar in the future, in particular with the advance of discrete/continuous optimization techniques (Lempitsky et al. 2008; Bleyer et al. 2010).

6 Conclusion

We have presented a collection of datasets for the evaluation of optical flow algorithms. These datasets are significantly more challenging and comprehensive than previous ones. We have also extended the set of evaluation measures and improved the evaluation methodology of Barron et al. (1994). The data and results are available at <http://vision.middlebury.edu/flow/>. Since the publication of our preliminary paper (Baker et al. 2007), a large number of authors have uploaded results to our online evaluation. The best results are a huge improvement over the algorithms in Baker et al. (2007) (Table 2). Our data and metrics are diverse, offering a number of insights into the choice of the most appropriate metrics and statistics (Sect. 5.2), the effect of the datatype on the performance of algorithms and the difficulty of the various forms of data (Sect. 5.2.4), the differences between flow errors and interpolation errors (Sect. 5.3), and the importance of the various components in an algorithm (Sect. 5.5). Of course, as newer papers continue to be published, e.g., Sun et al. (2010), which as we go to press (June 2010) is now the leading algorithm, our understanding of which factors contribute to good performance will continue to evolve.

Progress on our data has been so rapid that the performance on some of the sequences is already very good (Table 2). The main exceptions are *Grove*, *Teddy*, *Urban*, and perhaps *Schefflera*. As our statistical analysis shows, however, the correlation in performance across datasets is relatively low. This suggests that no *single* method is yet able to achieve strong performance across a wide variety of datatypes. We believe that such generality is a requirement for robust optical flow algorithms suited for real-world applications.

Any such dataset and evaluation has a limited lifespan and new and more challenging sequences should be collected. A natural question, then, is how such data is best collected. Of the various possible techniques—synthetic data (Barron et al. 1994; McCane et al. 2001), some form of hidden markers (Mova LLC 2004; Tappen et al. 2006; Ramnath et al. 2008), human annotation (Liu et al. 2008), interpolation data (Szeliski 1999), and modified stereo data (Scharstein and Szeliski 2003)—the authors believe that synthetic data is probably the best approach (although generating high-quality synthetic data is not as easy as it might

seem). Large motion discontinuities and fast motion of complex, fine structures appear to be more of a problem for current optical flow algorithms than non-rigid motion, complex illumination changes, and sensor noise. The level of difficulty is easier to control using synthetic data. Degradations such as sensor noise, etc., can also easily be added. The realism of synthetic sequences could also be improved further beyond the data in our evaluation.

Future datasets should also consider more challenging types of materials, illumination change, atmospheric effects, and transparency. Highly specular and transparent materials present not just a challenge for current algorithms, but also for quantitative evaluation. Defining the ground-truth flow and error metrics for these situations will require some care.

With any synthetic dataset, it is important to understand how representative it is of real data. Hence, the use of multiple types of data and an analysis of the correlation across them is critical. A diverse set of datatypes also reduces overfitting to any one type, while offering insights into the relative performance of the algorithms in different scenarios. On balance, however, we would recommend that any future studies contain a higher proportion of challenging, realistic synthetic data. Future studies should also extend the data to longer sequences than the 8-frame sequences that we collected.

Acknowledgements Many thanks to Brad Hiebert-Treuer and Alan Lim for their help in creating the fluorescent texture data sets. Michael Black and Stefan Roth were supported by NSF grants IIS-0535075 and IIS-0534858, and a gift from Intel Corporation. Daniel Scharstein was supported by NSF grant IIS-0413169. Aghiles Kheffache generously donated a software license for the 3Delight renderer for use on this project. Michael Black and JP Lewis thank Lance Williams for early discussions on synthetic flow databases and Doug Creel and Luca Fascone for discussions of rendering issues. Thanks to Sing Bing Kang, Simon Winder, and Larry Zitnick for providing implementations of various algorithms. Finally, thanks to all the authors who have used our data and uploaded results to our website.

Open Access This article is distributed under the terms of the Creative Commons Attribution Noncommercial License which permits any noncommercial use, distribution, and reproduction in any medium, provided the original author(s) and source are credited.

References

- Adiv, G. (1985). Determining three-dimensional motion and structure from optical flow generated by several moving objects. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 7(4), 384–401.
- Aggarwal, J., & Nandhakumar, N. (1988). On the computation of motion from sequences of images—a review. *Proceedings of the IEEE*, 76(8), 917–935.
- Anandan, P. (1989). A computational framework and an algorithm for the measurement of visual motion. *International Journal of Computer Vision*, 2(3), 283–310.

- Anandan, P., & Weiss, R. (1985). Introducing smoothness constraint in a matching approach for the computation of displacement fields. In *Proceedings of the DARPA image understanding workshop* (pp. 186–196).
- Baker, S., & Matthews, I. (2004). Lucas-Kanade 20 years on: a unifying framework. *International Journal of Computer Vision*, 46(3), 221–255.
- Baker, S., Scharstein, D., Lewis, J., Roth, S., Black, M., & Szeliski, R. (2007). A database and evaluation methodology for optical flow. In *Proceedings of the IEEE international conference on computer vision*.
- Barron, J., Fleet, D., & Beauchemin, S. (1994). Performance of optical flow techniques. *International Journal of Computer Vision*, 12(1), 43–77.
- Battiti, R., Amaldi, E., & Koch, C. (1991). Computing optical flow across multiple scales: an adaptive coarse-to-fine strategy. *International Journal of Computer Vision*, 6(2), 133–145.
- Beier, T., & Neely, S. (1992). Feature-based image metamorphosis. In *Annual conference series: Vol. 26(2). ACM computer graphics, SIGGRAPH* (pp. 35–42).
- Bergen, J., Anandan, P., Hanna, K., & Hingorani, R. (1992). Hierarchical model-based motion estimation. In *Proceedings of the European conference on computer vision* (pp. 237–252).
- Black, M., & Anandan, P. (1991). Robust dynamic motion estimation over time. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 296–302).
- 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.
- Black, M., & Jepson, A. (1996). Estimating optical flow in segmented images using variable-order parametric models with local deformations. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 18(10), 972–986.
- Blake, A., & Zisserman, A. (1987). *Visual reconstruction*. Cambridge: MIT Press.
- Bleyer, M., & Chambon, S. (2010). Does color really help in dense stereo matching? In *Proceedings of the international symposium 3D data processing, visualization and transmission*.
- Bleyer, M., Rother, C., & Kohli, P. (2010). Surface stereo with soft segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*.
- Boykov, Y., Veksler, O., & Zabih, R. (2001). Fast approximate energy minimization via graph cuts. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 23(11), 1222–1239.
- Brox, T., Bregler, C., & Malik, J. (2009). Large displacement optical flow. In *Proceedings of the IEEE conference on computer vision and pattern recognition*.
- Brox, T., Bruhn, A., Papenberg, N., & Weickert, J. (2004). High accuracy optical flow estimation based on a theory for warping. In *Proceedings of the European Conference on Computer Vision* (Vol. 4, 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.
- Bruhn, A., Weickert, J., Kohlberger, T., & Schnörr, C. (2006). A multigrid platform for real-time motion computation with discontinuity-preserving variational methods. *International Journal of Computer Vision*, 70(3), 257–277.
- Burt, P., Yen, C., & Xu, X. (1982). Local correlation measures for motion analysis: a comparative study. In *Proceedings of the IEEE conference on pattern recognition and image processing* (pp. 269–274).
- Burt, P., Yen, C., & Xu, X. (1983). Multi-resolution flow-through motion analysis. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 246–252).
- Cassisa, C., Simoens, S., & Prinet, V. (2009). Two-frame optical flow formulation in an unwarped multiresolution scheme. In *Proceedings of the Iberoamerican congress on pattern recognition* (pp. 790–797).
- Cooke, T. (2008). Two applications of graph-cuts to image processing. In *Proceedings of digital image computing: techniques and applications* (pp. 498–504).
- DNA Research (2008). 3Delight rendering software. <http://www.3delight.com/>.
- Enkelman, W. (1986). Investigations of multigrid algorithms for the estimation of optical flow fields in image sequences. In *Proceedings of the workshop on motion: representations and analysis* (pp. 81–87).
- Everingham, M., Van Gool, L., Williams, C., Winn, J., & Zisserman, A. (2009). The PASCAL visual object classes challenge 2009. <http://www.pascal-network.org/challenges/VOC/voc2009/workshop/index.html>
- Fei-Fei, L., Fergus, R., & Perona, P. (2006). One-shot learning of object categories. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 28(4), 594–611.
- Fleet, D., & Jepson, A. (1990). Computation of component image velocity from local phase information. *International Journal of Computer Vision*, 5(1), 77–104.
- Fuh, C., & Maragos, P. (1989). Region-based optical flow estimation. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 130–135).
- Georghiades, A., Belhumeur, P., & Kriegman, D. (2001). From few to many: illumination cone models for face recognition under variable lighting and pose. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 23(6), 643–660.
- Glazer, F., Reynolds, G., & Anandan, P. (1983). Scene matching by hierarchical correlation. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 432–441).
- Glocker, B., Paragios, N., Komodakis, N., Tziritas, G., & Navab, N. (2008). Optical flow estimation with uncertainties through dynamic MRFs. In *Proceedings of the IEEE conference on computer vision and pattern recognition*.
- Golland, P., & Bruckstein, A. (1997). Motion from color. *Computer Vision and Image Understanding*, 68(3), 346–362.
- Gross, R., Matthews, I., Cohn, J., Kanade, T., & Baker, S. (2008). Multi-PIE. In *Proceedings of the international conference on automatic face and gesture recognition*.
- Hanna, K. (1991). Direct multi-resolution estimation of ego-motion and structure from motion. In *Proceedings of the IEEE workshop on visual motion* (pp. 156–162).
- Haussecker, H., & Fleet, D. (2000). Computing optical flow with physical models of brightness variation. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (Vol. 2, pp. 760–767).
- Herbst, E., Seitz, S., & Baker, S. (2009). *Occlusion reasoning for temporal interpolation using optical flow*. Technical report UW-CSE-09-08-01, Department of Computer Science and Engineering University of Washington.
- Horn, B. (1986). *Robot vision*. Cambridge: MIT Press.
- Horn, B., & Schunck, B. (1981). Determining optical flow. *Artificial Intelligence*, 17, 185–203.
- Jepson, A., & Black, M. (1993). Mixture models for optical flow computation. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 760–761).
- Ju, S. (1998). *Estimating image motion in layers: the skin and bones model*. PhD thesis, Department of Computer Science, University of Toronto.
- Ju, S., Black, M., & Jepson, A. (1996). Skin and bones: multi-layer, locally affine, optical flow and regularization of transparency. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 307–314).

- Jung, H., Lee, K., & Lee, S. (2008). Toward global minimum through combined local minima. In *Proceedings of the European conference on computer vision* (Vol. 4, pp. 298–311).
- Landis, H. (2002). Production-ready global illumination. In L. Gritz (Ed.), *RenderMan in production: SIGGRAPH 2002 course 16* (pp. 87–100). New York: ACM.
- Le Besnerais, G., & Champagnat, F. (2005). Dense optical flow by iterative local window registration. In *Proceedings of the international conference on image processing* (Vol. 1, pp. 137–140).
- Lei, C., & Yang, Y. (2009). Optical flow estimation on coarse-to-fine region-trees using discrete optimization. In *Proceedings of the IEEE international conference on computer vision*.
- Lempitsky, V., Roth, S., & Rother, C. (2008). Fusion flow: discrete-continuous optimization for optical flow estimation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*.
- Levoy, M. (1988). Display of surfaces from volume data. *IEEE Computer Graphics and Applications*, 8(3), 29–37.
- Li, Y., & Huttenlocher, D. (2008). Learning for optical flow using stochastic optimization. In *Proceedings of the European conference on computer vision* (Vol. 2, pp. 373–391).
- Liu, C., Freeman, W., Adelson, E., & Weiss, Y. (2008). Human-assisted motion annotation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*.
- Liu, C., Yuen, J., Torralba, A., Sivic, J., & Freeman, W. (2008). SIFT flow: dense correspondence across difference scenes. In *Proceedings of the European conference on computer vision* (Vol. 3, pp. 28–42).
- Lucas, B., & Kanade, T. (1981). An iterative image registration technique with an application in stereo vision. In *Proceedings of the international joint conference on artificial intelligence* (pp. 674–679).
- Mahajan, D., Huang, F., Matusik, W., Ramamoorthi, R., & Belhumeur, P. (2009). Moving gradients: a path-based method for plausible image interpolation. In *Annual conference series. ACM computer graphics, SIGGRAPH*.
- Markandey, V., & Flinchbaugh, B. (1990). Multispectral constraints for optical flow computation. In *Proceedings of the IEEE international conference on computer vision* (pp. 38–41).
- McCane, B., Novins, K., Crannitch, D., & Galvin, B. (2001). On benchmarking optical flow. *Computer Vision and Image Understanding*, 84(1), 126–143.
- Mitiche, A., & Boutheny, P. (1996). Computation and analysis of image motion: a synopsis of current problems and methods. *International Journal of Computer Vision*, 19(1), 29–55.
- Mova LLC (2004). Contour reality capture. <http://www.mova.com/>.
- Murray, D., & Buxton, B. (1987). Scene segmentation from visual motion using global optimization. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 9(2), 220–228.
- Nagel, H.-H., & Enkelmann, W. (1986). An investigation of smoothness constraints for the estimation of displacement vector fields from image sequences. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 8(5), 565–593.
- 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*, 20(9), 961–979.
- Nir, T., Bruckstein, A., & Kimmel, R. (2008). Over-parameterized variational optical flow. *International Journal of Computer Vision*, 76(2), 205–216.
- Ohta, N. (1989). Optical flow detection by color images. In *International conference on image processing* (pp. 801–805).
- Otte, M., & Nagel, H.-H. (1994). Optical flow estimation: advances and comparisons. In *Proceedings of the European conference on computer vision* (pp. 51–60).
- Philips, P., Scruggs, W., O'Toole, A., Flynn, P., Bowyer, K., Schott, C., & Sharpe, M. (2005). Overview of the face recognition grand challenge. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (Vol. 1, pp. 947–954).
- Pock, T., Pock, M., & Bischof, H. (2007). Algorithmic differentiation: application to variational problems in computer vision. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 29(7), 1180–1193.
- Pratt, W. (1974). Correlation techniques of image registration. *IEEE Transactions on Aerospace and Electronic Systems, AES-10*, 353–358.
- Ramnath, K., Baker, S., Matthews, I., & Ramanan, D. (2008). Increasing the density of active appearance models. In *Proceedings of the IEEE conference on computer vision and pattern recognition*.
- Rannacher, J. (2009). *Realtime 3D motion estimation on graphics hardware*. Undergraduate thesis, Heidelberg University.
- Ren, X. (2008). Local grouping for optical flow. In *Proceedings of the IEEE conference on computer vision and pattern recognition*.
- Roth, S., & Black, M. (2007). On the spatial statistics of optical flow. *International Journal of Computer Vision*, 74(1), 33–50.
- Scharstein, D., & Szeliski, R. (2002). A taxonomy and evaluation of dense two-frame stereo correspondence algorithms. *International Journal of Computer Vision*, 47(13), 7–42.
- Scharstein, D., & Szeliski, R. (2003). High-accuracy stereo depth maps using structured light. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 195–202).
- Seitz, S., & Baker, S. (2009). Filter flow. In *Proceedings of the IEEE international conference on computer vision*.
- Seitz, S., Curless, B., Diebel, J., Scharstein, D., & Szeliski, R. (2006). A comparison and evaluation of multi-view stereo reconstruction algorithms. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (Vol. 1, pp. 519–526).
- Shade, J., Gortler, S., He, L.-W., & Szeliski, R. (1998). Layered depth images. In *Annual conference series. ACM computer graphics, SIGGRAPH* (pp. 231–242).
- Shizawa, M., & Mase, K. (1991). A unified computational theory for motion transparency and motion boundaries based on eigenenergy analysis. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 289–295).
- Sim, T., Baker, S., & Bsas, M. (2003). The CMU pose, illumination, and expression database. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 25(12), 1615–1618.
- Stiller, C., & Konrad, J. (1999). Estimating motion in image sequences: a tutorial on modeling and computation of 2D motion. *IEEE Signal Processing Magazine*, 16(4), 70–91.
- Sun, C. (1999). Fast optical flow using cross correlation and shortest-path techniques. In *Proceedings of digital image computing: techniques and applications* (pp. 143–148).
- Sun, J., Shum, H.-Y., & Zheng, N. (2003). Stereo matching using belief propagation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 25(7), 787–800.
- Sun, D., Roth, S., Lewis, J., & Black, M. (2008). Learning optical flow. In *Proceedings of the European conference on computer vision* (Vol. 3, pp. 83–97).
- Sun, D., Roth, S., & Black, M. (2010). Secrets of optical flow estimation and their principles. In *Proceedings of the IEEE conference on computer vision and pattern recognition*.
- Szeliski, R. (1999). Prediction error as a quality metric for motion and stereo. In *Proceedings of the IEEE international conference on computer vision* (pp. 781–788).
- Tappen, M., Adelson, E., & Freeman, W. (2006). Estimating intrinsic component images using non-linear regression. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (Vol. 2, pp. 1992–1999).
- Trobin, W., Pock, T., Cremers, D., & Bischof, H. (2008). Continuous energy minimization via repeated binary fusion. In *Proceedings*

- of the European conference on computer vision* (Vol. 4, pp. 677–690).
- Trobin, W., Pock, T., Cremers, D., & Bischof, H. (2008). An unbiased second-order prior for high-accuracy motion estimation. In *Proceedings of pattern recognition, DAGM* (pp. 396–405).
- Valgaerts, L., Bruhn, A., & Weickert, J. (2008). A variational model for the joint recovery of the fundamental matrix and the optical flow. In *Proceedings of pattern recognition, DAGM* (pp. 314–324).
- Vedula, S., Baker, S., Rander, P., Collins, R., & Kanade, T. (2005). Three-dimensional scene flow. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 27(3), 475–480.
- Wang, J., & Adelson, E. (1993). Layered representation for motion analysis. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 361–366).
- Wedel, A., Pock, T., Braun, J., Franke, U., & Cremers, D. (2008). Duality TV-L1 flow with fundamental matrix prior. In *Proceedings of image and vision computing*, New Zealand.
- Wedel, A., Pock, T., Zach, C., Cremers, D., & Bischof, H. (2008). An improved algorithm for TV-L1 optical flow. In *Proceedings of the Dagstuhl motion workshop*.
- Wedel, A., Cremers, D., Pock, T., & Bischof, H. (2009). Structure-and motion-adaptive regularization for high accuracy optic flow. In *Proceedings of the IEEE international conference on computer vision*.
- Weiss, Y. (1997). Smoothness in layers: motion segmentation using nonparametric mixture estimation. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 520–526).
- Werlberger, M., Trobin, W., Pock, T., Bischof, H., Wedel, A., & Cremers, D. (2009). Anisotropic Huber-L1 optical flow. In *Proceedings of the British machine vision conference*.
- Xu, L., Chen, J., & Jia, J. (2008). A segmentation based variational model for accurate optical flow estimation. In *Proceedings of the European conference on computer vision* (Vol. 1, pp. 671–684).
- Zimmer, H., Bruhn, A., Weickert, J., Valgaerts, L., Salgado, A., Rosenhahn, B., & Seidel, H.-P. (2009). Complementary optic flow. In *Proceedings of seventh international workshop on energy minimization methods in computer vision and pattern recognition*.
- Zitnick, C., Kang, S., Uyttendaele, M., Winder, S., & Szeliski, R. (2004). High-quality video view interpolation using a layered representation. In *Annual conference series: Vol. 23(2). ACM computer graphics, SIGGRAPH* (pp. 600–608).