

# Learning a Confidence Measure for Optical Flow

Oisín Mac Aodha, *Student Member, IEEE*, Ahmad Humayun, *Student Member, IEEE*, Marc Pollefeys, *Fellow, IEEE*, and Gabriel J. Brostow, *Member, IEEE*

**Abstract**—We present a supervised learning-based method to estimate a per-pixel confidence for optical flow vectors. Regions of low texture and pixels close to occlusion boundaries are known to be difficult for optical flow algorithms. Using a spatiotemporal feature vector, we estimate if a flow algorithm is likely to fail in a given region. Our method is not restricted to any specific class of flow algorithm and does not make any scene specific assumptions. By automatically learning this confidence, we can combine the output of several computed flow fields from different algorithms to select the best performing algorithm per pixel. Our optical flow confidence measure allows one to achieve better overall results by discarding the most troublesome pixels. We illustrate the effectiveness of our method on four different optical flow algorithms over a variety of real and synthetic sequences. For algorithm selection, we achieve the top overall results on a large test set, and at times even surpass the results of the best algorithm among the candidates.

**Index Terms**—Optical flow, confidence measure, Random Forest, synthetic data, algorithm selection

## 1 INTRODUCTION

BENCHMARKING datasets such as the Middlebury Optical Flow Evaluation Table [1] have motivated improvements in the accuracy of optical flow algorithms. These evaluations, while also useful for highlighting areas of future research, can still leave practitioners uncertain about how to capitalize on the rankings. It is difficult for most nonexperts to assess how suitable a particular algorithm will be, given their data. The expense and difficulty of obtaining ground truth for real-world scenes when evaluating algorithm/scene pairings is enormous. This leaves practitioners trying to choose which among the very few image pairs is most like their test video at hand. To a limited extent, each algorithm can be used to self-assess its own performance. Algorithms that seek to optimize a nonconvex energy term at test time know only that a local optimum has been reached once they have converged. This energy state is often interpreted as a confidence, but it is not directly comparable between several different algorithms due to different energy terms or priors being utilized. In our previous work, we introduced a supervised learning-based confidence measure for optical flow that

gives us a probabilistic estimate of confidence [2]. We do not rely on any scene assumptions and our confidence can be computed for any type of flow algorithm. In this paper, we improve and carefully measure the accuracy of our confidence measure. We also propose a meta-algorithm that automatically chooses the most appropriate algorithm for the situation.

We define confidence,  $\psi$ , for each flow vector as the probability of that flow being below some specified error threshold  $\epsilon_{epe}^s$ , where  $\epsilon_{epe}^s$  is the amount of end point error (EPE) acceptable to the user. Confidence measures for optical flow have been explored in the past. However, they have typically been algorithm-type specific [4] or have made simplifying assumptions about the statistics of local flow [5]. We seek out the correlation between good performance by a constituent algorithm and specific local situations that can be discerned statistically from the image sequence. Fig. 1 illustrates a typical confidence image from our algorithm.

The semantic segmentation community has been developing successful techniques to find correlations between object classes and appearance (e.g., [6] and [7]). Using similar intuition, we learn the relationship between spatiotemporal image features and algorithm success. We attempt to predict the best algorithm locally, given a set of candidate flow algorithms, where the “best” algorithm is the one that is predicted to yield the best accuracy. We assume that implementations of all the algorithms under consideration are available. Recognizing that most flow algorithms may be ported to leverage GPU processing, we accept the fixed cost of running all of them on a given sequence as acceptable in pursuing the best overall accuracy. We extend our previous work [2] and make the following contributions:

- O. Mac Aodha and G.J. Brostow are with the Department of Computer Science, University College London, Gower Street, London WC1E 6BT, United Kingdom. E-mail: {o.macaodha, G.Brostow}@cs.ucl.ac.uk.
- A. Humayun is with the School of Interactive Computing, Georgia Institute of Technology, 304N CCB, Atlanta, GA 30332. E-mail: ahumayun@cc.gatech.edu.
- M. Pollefeys is with the Computer Vision and Geometry Laboratory, Department of Computer Science, ETH Zurich, CNB G105, Universitätsstrasse 6, CH-8092 Zurich, Switzerland. E-mail: marc.pollefeys@inf.ethz.ch.

Manuscript received 23 Feb. 2012; revised 2 July 2012; accepted 25 July 2012; published online 2 Aug. 2012.

Recommended for acceptance by C. Lampert.

For information on obtaining reprints of this article, please send e-mail to: [tpami@computer.org](mailto:tpami@computer.org), and reference IEEECS Log Number TPAMI-2012-02-0135.

Digital Object Identifier no. 10.1109/TPAMI.2012.171.

- a thorough evaluation of our optical flow confidence measure on new flow algorithms and several new sequences,
- comparison to other baseline confidence measures,

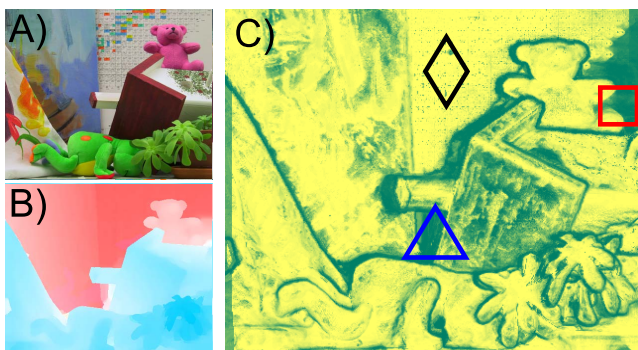


Fig. 1. Optical flow confidence. A) Input image, one of two. B) Computed flow field using [3]. C) Confidence image: Green indicates low confidence, while yellow is high. Our algorithm correctly identifies confidence for situations such as  $\triangle$  motion discontinuities,  $\diamond$  high and  $\square$  low texture.

- separate confidence in the  $X$  and  $Y$  directions,
- improved accuracy for optical flow by automatically selecting among known constituent algorithms,
- an improved system for easily producing synthetic ground-truth optical flow data for scenes with moving objects.

Experiments show our confidence measure outperforms other general purpose measures. Additionally, we automatically combine the output of multiple different algorithms, which gives better results than any individual algorithm.

## 2 RELATED WORK

We examine the relevant work in optical flow *confidence estimation*. For an overview of traditional optical flow approaches see [8], and [1] or [3] for more current techniques. We also review work related to algorithm selection, defined as finding the algorithm from a candidate set that produces the most accurate result for a given task-algorithm combination.

### 2.1 Confidence Estimation

Early confidence measures for optical flow were only concerned with intensity information. Simoncelli et al. [9] proposed a method based on the gradient of the intensity in a window about the patch. The justification is that one would expect computed flow to be accurate in areas of high gradient, e.g., high texture regions and image corners. Their method does not just return a single confidence estimate for each vector, but a 2D distribution which they use to represent uncertainty. Anandan [10] also expresses confidence as a 2D measure of the curvature in the sum of squared differences surface computed during candidate matching. This choice of 2D confidence is that it can represent the certainty of the flow in a particular direction (both  $x$  and  $y$ ). Uras et al. [11] look at the spatial Hessian matrix of the local intensity patch. Jähne et al. [12] present several methods based on an eigenvalue decomposition of the 3D structure tensor. Some of their measures look at the temporal gradient, but do not take the computed flow field into account. In effect, these measures attempt to predict how difficult it will be to determine flow for a particular image pair by analyzing their spatial and temporal gradients. Our approach differs in that it learns a mapping

between flow algorithm success and the spatiotemporal image data.

Algorithm-specific confidence estimation techniques also exist. Kybic and Nieuwenhuis [4] describe a method which works for optical flow algorithms that minimize spatially decomposable variational image similarity terms, such as [13], [14]. Their bootstrap resampling approach must compute the flow field over multiple iterations (10 in their paper), while at each iteration the input data are perturbed and the variability of the result is measured. As noted by the authors, their algorithm may succeed in detecting the variance in the error but not the bias. Bruhn and Weickert [15] propose a confidence measure for variational optical flow methods where confidence is inversely proportional to the local energy of the objective being minimized. For other examples of algorithm specific methods which do not generalize across optical flow algorithms, see [4], [16], [17].

Kondermann et al. [18] propose a PCA-based method where confidence is defined as how well a learned linear subspace approximates the test flow vector. In follow up work [5], they learn a probabilistic model of the flow field in local windows from training data. These flow vectors are then modeled as a multivariate Gaussian distribution and a confidence measure is proposed based on statistical test theory. These two approaches are most closely related to ours in that they attempt to learn a model from training data, but differ in the fact that they rely on strong assumptions regarding local smoothness. In our previous work [2], which we build on here, we used a supervised learning approach to estimate confidence for both interest point descriptors and optical flow. Our method seeks to learn where each flow algorithm will succeed or fail based on analyzing a feature vector computed from the image pair. We combine multiple feature types such as temporal, texture, distance from image edges, and others, to estimate the confidence in a given flow algorithm's success. This confidence was also employed in a state-of-the-art occlusion detector [19]. Gehrig and Scharwächter describe a real-time system which combines several different cues (with one variational flow specific feature) to estimate confidence by classifying pixels into discrete error classes based on estimated flow error [20].

Attempts have been made to evaluate the performance of different confidence measures. Bainbridge-Smith and Lane [21] compare several spatial derivative-based confidence measures on a limited set of data. Kybic and Nieuwenhuis [4] provide a thorough comparison of their work against others but only for one optical flow algorithm.

Other areas have witnessed attempts to learn a confidence measure. For depth images captured using a Time Of Flight camera, Reynolds et al. [22] proposed a supervised learning method in the spirit of this work, which classifies the depth error returned by the camera. Using a learning-based approach, Li et al. [23] sought to learn a ranking function which sorts interest points according to their stability.

### 2.2 Algorithm Selection

In addition to estimating confidence for a particular flow algorithm, our supervised learning approach also allows us to combine the output of several different flow algorithms to choose the best flow at each pixel. Here, we review related

work in combining different “experts” with specific emphasis on methods for combining optical flow algorithms.

Raykar et al. [24] proposed a model to deal with the scenario in supervised learning where multiple annotators (or experts) exist, but each of them is slightly wrong. In their scenario, *one* expert is assumed to always be better than all the rest, and the task consists of finding that expert. The technique is an improvement over following the majority vote when some experts are better than others. Our problem formulation is different, however, because we cannot assume that one expert is consistently better or worse, independent of the image data being considered.

Learned algorithm selection is shown by Yong et al. [25] for the specific task of image segmentation. They used an SVM for learning and performed their experiments on 1,000 synthetic images of 2D squares, circles, etc., with additive noise, demonstrating what is actually online learning for algorithm selection. Working with 14 constituent real-time tracking algorithms, Stenger et al. [26] developed a framework that learned the expected error of each algorithm, given its confidence values. Then, during testing, the best performing pairs of algorithms could be cascaded or run in parallel to track a hand or head. This approach is very flexible for situations where one task is being accomplished at a time. Alt et al. [27] describe a supervised learning approach for assessing which planar patches will be difficult for tracking. Using this preselection of reliable templates, they report an improved detection rate for an existing tracking-by-detection system. Peng and Veksler attempt to automatically estimate the best parameters for interactive segmentation [28]. They train a classifier on image features computed from training data and, during testing, attempt to choose the best set of parameters (where a parameter set could be viewed as an algorithm) to segment the given scene.

Muja and Lowe [29] have presented a unique approach to algorithm selection that is quite valuable in the context of feature matching and beyond. Like us, they argue that algorithm suitability is data dependent. Their system searches a parameter space, where the algorithm itself is just one of the parameters, to find an appropriate approximate nearest neighbor strategy (algorithm and settings). The automatically determined strategy is based on the target data itself, such as a database of SIFT descriptors [30], and desired preferences for optimizing lookup speeds versus memory. There, the training data are the same as the test data, so their optimization is deterministic, while our algorithm suitability must be learned so we can predict which segments are suited for which strategy just by looking at each video.

Of the existing approaches to computing optical flow, the iterative FusionFlow [31] is still very different technically, but the closest to our approach in terms of its philosophy. They compute a discrete optimization on continuous-valued flow-fields (with another continuous optimization “clean up”) by performing a minimal cut on an extended graph. The extended graph consists of auxiliary binary-valued labels to represent either accepting a newly proposed flow vector at that location or keeping the current flow estimate. The similarity to our work is that in each such iteration of FusionFlow, the new proposed solution could be viewed as a competing strategy or algorithm, offering a

potentially lower energy than the current estimate, at least in some spatial neighborhood. FusionFlow is quite flexible and could potentially be modernized with more competitive starting proposals than the 200+ based on Lucas and Kanade [32] and Horn and Schunk [33], but the authors indicate that because of their energy function, the computed minimum eventually gives a score extremely close to the energy of the ground-truth solution.

A thorough understanding of existing energy functions allowed Bruhn et al. [34] to formulate a new Combined Local-Global (CLG) method, aptly named “Lucas/Kanade Meets Horn/Schunck.” Their new 2D energy term (and its 3D variant) combined the local robustness to noise offered by algorithms such as Lucas and Kanade [32] with the regularized smoothness and dense flow of global algorithms such as Horn and Schunk [33]. They compute a confidence criterion based on this new energy term, and demonstrate that it is partly correlated with actual accuracy. The challenge they describe has been one of our driving motivations, namely, that one has few if any reliable confidence measures beyond the chosen energy function itself. That problem is compounded when comparing multiple algorithms with different energy-minimization objectives.

The nonparametric FRAME model of Zhu et al. [35] optimized its texture synthesis by picking out filters from a filter bank whose responses are correlated with neighborhoods in the training image. That approach is very flexible, adaptively using potentially many filters, including nonlinear ones which filter large subimages. Since then, Roth and Black’s Fields of Experts (FoE) [36] has gained a following by augmenting FRAME, extending Markov random fields with the capability to *learn* filters that model local field potentials. The completely data-driven nature of FoE is very attractive, and Woodford et al. [37] showed a method that trains with  $5 \times 5$  cliques in a comparatively short time. Roth and Black have further demonstrated FoE for the purpose of modeling an optical flow prior [38]. In [38], they used range images of known scenes with separately obtained real camera motions to learn a model of motion fields which are different from optical flow. Here, they still had to manually monitor convergence of the learning, but in testing, demonstrated superior results using these spatial statistics as priors for the aforementioned 2D Bruhn et al. [34] flow algorithm. FoE’s expert functions are less flexible than the FRAME model by design: They can be nonlinear, but need to be continuous, and the log of an expert has to be differentiable with respect to both the expert’s parameters and the (necessarily) linear filter responses.

Sun et al. [39] adapted their spatial FoE model of optical flow, learning a relationship between image and flow boundaries, this time with a parameterization of spatio-temporal brightness inconstancy. The steered model of flow and the generalized data term are learned on the painstakingly prepared ground-truth flow data of Baker et al. [1]. In our experiments, we too train on similar data and also have no need for sequence-specific parameter tuning, and we achieve better scores simply by virtue of leveraging multiple black-box algorithms that are effective in their own right.

An important result of the FoE line of research is the finding that, with careful optimization procedures, a good generalist algorithm's priors about local responses to linear filters should be learned from representative training data. Different low-dimensional "experts" in this setting are not unique algorithms, but are instead measures being combined to model high-dimensional probability distributions of parameterized statistics. Our goal is much simpler, nonparametric, and complementary: to establish the *discriminability* between visual situations given competing strategies or algorithms, in this case, for computing optical flow. For example, the algorithms with FoE-based priors trained with different sized cliques ( $5 \times 5$  for [38],  $9 \times 9$  for [39]) could be evaluated as different strategies in our framework.

### 3 LEARNING ALGORITHM

Given a dense optical flow field  $F$ , computed from an image pair  $I_1$  and  $I_2$ , we wish to estimate a confidence value  $\psi^i \in [0, 1]$  for each flow vector  $\mathbf{f}_i = (u_i, v_i)$ . One option would be to pose this as a regression task and attempt to estimate the true error value  $\epsilon_{epe}^*$  for each flow vector, where  $\epsilon_{epe}^*$  is the End Point Error, i.e., the distance measured in pixels between the computed flow vector and the ground truth. Instead, we attempt to solve the comparatively easier problem of determining if the proposed flow vector  $\mathbf{f}_i$  is reliable or not at a specific error threshold  $\epsilon_{epe}^s$ . Unlike other methods, this has the advantage of allowing the user to specify a lower limit on accuracy. For example, in some applications, it is beneficial to have more pixels, even with coarser flow estimates, e.g., [40]. We pose confidence estimation as a standard binary supervised learning problem of the form

$$\mathcal{D} = \{(\mathbf{x}_i, c_i) | \mathbf{x}_i \in \mathbb{R}^d, c_i \in \{0, 1\}\}_{i=1}^n, \quad (1)$$

with  $n$  being the number of training examples,  $d$  the dimensionality of the feature vector  $\mathbf{x}_i$  computed from the images and flow field, and  $c_i$  the label. In training, a flow vector  $\mathbf{f}_i$  gets a label of 1 if its EPE,  $\epsilon_{epe}^i$ , is less than the desired threshold  $\epsilon_{epe}^s$ ; otherwise it is set to 0:

$$c_i = \begin{cases} 1 & \epsilon_{epe}^i \leq \epsilon_{epe}^s \\ 0 & \epsilon_{epe}^i > \epsilon_{epe}^s \end{cases} \quad (2)$$

At test time, the probability associated with the class label  $c_i$  is taken to be our confidence  $\psi^i$ .

The applicability of most flow algorithms is situation specific, and we wish to classify those situations automatically. Using a similar approach, we seek to learn the mapping between a feature vector and a class label which represents the different possible algorithms. In this scenario, algorithm selection is posed as a multiclass supervised learning problem:

$$\mathcal{D} = \{(\mathbf{x}_i, c_i) | \mathbf{x}_i \in \mathbb{R}^d, c_i \in \mathbb{Z}^K\}_{i=1}^n, \quad (3)$$

with the same notation as (1), but now  $c_i$  is the algorithm with the lowest EPE and  $K$  is the number of possible competing algorithms.

Our single classifier is taking the place of the multiple algorithm-specific energy terms or confidence measures. Being probabilistic, the posteriors of different classifiers

can be compared to each other. Task accuracy should be improved if each part of an image sequence is handled by the most suitable of  $K$  algorithms. The proposed approach is most appropriate in situations where either no good single algorithm exists or where a generalist algorithm makes mistakes in places that some specialist algorithm does not.

### 3.1 Choice of Algorithm

For our classifier, we have selected the Random Forests algorithm developed by Breiman [41]. Random Forests is an ensemble of decision trees which averages the predictions of the trees to assign the class labels. It makes use of bagging to uniformly sample (with replacement) subsets from the dataset to train the decision trees. It can also use the remaining data to estimate the error for that particular tree. During training, each node selects from a random set of tests the one that best splits that data. A Random Forest has the advantage of being fast to train and test even on large amounts of data; it is multiclass, robust to noise, inherently parallelizable, can handle large datasets, and it also estimates the importance of the input variables. See [42] for a detailed overview of classification and regression forests. We also experimented with Boosted Trees and SVMs and noted slightly worse performance with an increase in training time.

## 4 FEATURES

Given an image pair  $I_1$  and  $I_2$  (where  $I = f(x, y)$  is a grayscale image), we wish to construct a feature representation for each pixel in the first image,  $\mathbf{x}_i$ , which is indicative of the success and failure cases of optical flow algorithms. We use a similar feature representation to [2], with the addition of some new features from [19]. This feature set, while certainly not exhaustive, combines single image, temporal, and scale space features.

### 4.1 Appearance

Highly textured regions provide little challenge for modern optical flow algorithms. By taking the gradient magnitude of the image, it is possible to measure the level of "texturedness" of a region:

$$g(x, y, z) = \|\nabla I_1(x, y, z)\|, \quad (4)$$

where  $x$  and  $y$  are the pixel location in  $I_1$ , and  $z$  is the level in the image pyramid. Additionally, the distance transform is calculated on Canny edge detected images:

$$d(x, y, z) = \text{disTrans}(\|\nabla I_1(x, y, z)\| > \tau_{ed}). \quad (5)$$

The intuition is that image edges may co-occur with motion boundaries, and the higher the distance from them, the lower the chance of occlusion. We also use the learned  $P_b$  edge detector of [43], which produces edge maps that often correlate with object edges:

$$pb(x, y, z) = \text{disTrans}(P_b[I_1(x, y, z)] > \tau_{pb}). \quad (6)$$

Other texture-based features, such as convolution with filter banks, were tested to capture other neighborhood information, but did not show increased performance.

## 4.2 Temporal

Flow algorithms tend to break down at motion discontinuities. Identifying these regions can be a cue for improving flow accuracy. Techniques such as image differencing can potentially locate these regions, but we found that a more robust approach is to take the derivative of the proposed flow fields. This is done by computing the median of the different candidate algorithms' flow and then calculating the gradient magnitude in the  $x$  and  $y$  directions, respectively:

$$t_x(x, y, z) = \|\nabla \dot{u}\|, \quad t_y(x, y, z) = \|\nabla \dot{v}\|. \quad (7)$$

## 4.3 Photo Constancy

Another indicator of optical flow quality is to measure the photoconstancy residual. For a given pixel, this is achieved by subtracting the intensity in  $I_2$  at  $x, y$  advected with the predicted flow  $u, v$  from the intensity in  $I_1$  at  $x, y$ . Due to the discrete nature of image space, we bicubically interpolate the intensity values in the second image. The residual error, measured in intensity, is calculated independently for each of the  $K$  candidate flow algorithms, so

$$r(x, y, k) = |I_1(x, y) - \text{bicubic}(I_2(x + u^k, y + v^k))|. \quad (8)$$

In the scenario where the optical flow vector projects the pixel outside the bounds of  $I_2$ , we assign a constant penalty.

## 4.4 Scale

Most effective approaches to optical flow estimation utilize scale space to compute flow for big motions. With this in mind, all of these features, with the exception of the residual error, are calculated on an image pyramid with  $Z = \{1, \dots, l\}$  levels and a rescaling factor of  $s$ .

These individual features are combined to create the full feature vector  $\mathbf{x}_i$ , computed for each of the pixels in  $I_1$  as

$$\mathbf{x}_i = \{g(x, y, Z), d(x, y, Z), t_x(x, y, Z), t_y(x, y, Z), pb(x, y, Z), r(x, y, \{1, \dots, k\})\}. \quad (9)$$

## 5 TRAINING DATA

Several techniques have been proposed to generate ground-truth optical flow data from real-image sequences. The popular Middlebury optical flow dataset approximated flow by painting a scene with hidden fluorescent texture and imaging it under UV illumination [1]. The ground-truth flow is then computed by tracking small windows in the high-resolution UV images and performing a brute-force search in the next frame. The high-resolution flow field is then downsampled to produce the final ground truth. This technique, while successful, is extremely time consuming and limited in the types of scenes that can be captured (restricted to lab environments). Additionally, the ambiguity in matching the image patches can result in incorrect flow and inaccurate labeling of occlusion regions. Human assistance has been used to explicitly annotate motion boundaries in scenes [44]. However, these approaches remain inaccurate and not scalable for producing large amounts of reliable ground-truth data.

Synthetically generated data offers an attractive method for automatically creating large amounts of accurate training data. This type of approach has been shown to be successful



Fig. 2. Ground-truth optical flow data. The top row depicts some example images from our system. Below is the ground-truth flow between successive frames. The flow field color is coded using the same format as [1]. Black values in the flow image indicate areas of occlusion between the two frames.

in applications such as human body pose estimation [45] and depth superresolution [46]. Synthetically generated sequences have been used as an alternative to natural images for optical flow evaluation since the introduction of the famous Yosemite sequence by Barron et al. [8]. Until now, the limiting factor in their use has been the inability to easily generate realistic sequences. As a result, practitioners have focused on “toy” datasets with unrealistic geometry and lighting [47], [48]. Using realistic texture, global illumination techniques, and by modeling complex geometry, it is now possible to generate realistic sequences with consumer 3D computer graphics packages [49]. Attempts have been made to assess whether synthetic data produces the same error distribution as real data [50].

In our previous work, we presented a system which allowed the user to generate ground-truth optical flow for a given synthetic image pair of a static scene [2]. Our expanded system allows us to generate ground-truth flow for arbitrary scenes with rigid moving objects and camera motion. Examples of our training data are shown in Fig. 2. The system works by casting a ray from the camera center in the first image through the image plane and into the scene until it intersects an object. Then, this point is projected back into the second camera (respecting occlusions in the scene and both camera and object motion), and the optical flow is calculated from the position difference with respect to the first image plane. An advantage of the system is that the texture and lighting of the scene is independent of the geometry. This creates the possibility for re-rendering the same scene using different illumination and textures, without altering the ground truth. As the system calculates intersections between projected rays and scene objects, occlusions are noted and therefore not erroneously labeled with incorrect flow (black regions in Fig. 2). To generate large amounts of data, we simulate rigid body dynamics on the scene objects with gravitational and force fields. We then randomly texture the objects from a library of high-resolution texture maps.

For our training set, we generated 20 image pairs. These scenes exhibit three different motion categories, small, medium, and large (determined by the median flow vectors). Seven of the scenes feature moving objects, and 14 have a moving camera. Camera motions include panning left and right, and rotation about the focal object. While certainly not an exhaustive set, these scenes are an attempt to cover a subspace of plausible scene motions. All training scenes, with descriptions of scene motion and texture content, with the code to create additional ground-truth data, are available on our project website.

## 5.1 Training Data Selection

Due to the potentially unlimited training data available, it is necessary to perform some selection on the examples used. In the algorithm selection case, training data can be quite redundant, as different algorithms can give the correct (or very close to correct) flow for a given pixel. Also, large portions of scenes can have very similar regions of flow (e.g., planar surfaces), offering little additional information. To overcome these problems, we preselect a subset of the available data on which to train. We only train on examples where the end point error between the best performing algorithm and the second best for a particular pixel is greater than a threshold. We set this threshold to a value of 0.3 pixels (based on the median difference between the best and worst algorithms for the whole dataset), which maximizes the number of training points where the constituent algorithms differ most. We also ensure that we have an equal amount of training data for each of the algorithms. For an experimental analysis of the effects of varying the amount of training data, see Section 7.2.1. In the case of confidence estimation, we select subsets at random from the training data (equally sampling from each scene). Samples which fall below  $\epsilon_{epe}^s$  are labeled as class 1 (acceptable error) and samples above are set to class 0 (too large an error). This reduces the amount of training data but also allows the selection of examples which are most discriminative.

## 6 COMPETING METHODS

We also compare our confidence measure against several competing methods. Like our results, each of these confidence measures is computed per pixel  $i$ . While additional confidence measures exist (e.g., [4]), we only consider those that are generally applicable to any type of flow algorithm.

The first and most basic measure attempts to characterize pixels of low texture because optical flow algorithms without any form of spatial regularization typically break down in these regions [16]. Here, confidence is related to the gradient magnitude of intensity in the first image:

$$\psi_{grad}^i = \|\nabla I_1\|. \quad (10)$$

The next set of confidence measures is based on properties of the 3D structure tensor [12]. The structure tensor  $\mathbf{J}$  is a 3D symmetric matrix of partial derivatives, computed from the spatiotemporal image sequence. Unlike the previous measure, these confidence measures use both images in the sequence to construct the structure tensor, though they still do not use any information specific to the flow computed. The structure tensor is computed for each pixel and has the form:

$$\mathbf{J}^i = \begin{bmatrix} \tilde{I}_{xx}^i & \tilde{I}_{xy}^i & \tilde{I}_{xt}^i \\ \tilde{I}_{xy}^i & \tilde{I}_{yy}^i & \tilde{I}_{yt}^i \\ \tilde{I}_{xt}^i & \tilde{I}_{yt}^i & \tilde{I}_{tt}^i \end{bmatrix}, \quad (11)$$

where  $\tilde{I}_{pq}^i$  is the smoothed<sup>1</sup> product of the partial derivatives in the  $p$  and  $q$  direction at pixel  $i$ . The derivatives are

1. In our experiments, smoothing is performed by convolving the derivatives with a  $7 \times 7$  Gaussian kernel with standard deviation 2.

approximated using finite differences in the  $x$ ,  $y$ , and  $t$  ( $I_1 \rightarrow I_2$ ) dimensions. An eigenvalue decomposition is then performed on this matrix, and the resulting eigenvalues ( $\lambda_1, \lambda_2, \lambda_3$ ) are used to compute the confidence. The eigenvalues are sorted into descending order, where  $\lambda_1 \geq \lambda_2 \geq \lambda_3 \geq 0$ .

The first structure tensor-based measure is the total coherency measure. It seeks to estimate the overall certainty of displacement:

$$\psi_{strTc}^i = \left( \frac{\lambda_1 - \lambda_3}{\lambda_1 + \lambda_3} \right)^2. \quad (12)$$

The spatial coherency measure seeks to detect the aperture problem, so

$$\psi_{strCs}^i = \left( \frac{\lambda_1 - \lambda_2}{\lambda_1 + \lambda_2} \right)^2. \quad (13)$$

The corner measure is computed as the difference between the previous two, so

$$\psi_{strCc}^i = \left( \frac{\lambda_1 - \lambda_3}{\lambda_1 + \lambda_3} \right)^2 - \left( \frac{\lambda_1 - \lambda_2}{\lambda_1 + \lambda_2} \right)^2. \quad (14)$$

The size of the smallest eigenvalue,  $\lambda_3$ , is correlated with homogeneous regions [4]:

$$\psi_{strEv3}^i = \lambda_3. \quad (15)$$

The previous structure tensor-based measures are agnostic to the computed flow fields. Kondermann et al. [5] describe a statistical test-based method that is trained on local examples of ground-truth optical flow. Unlike the previous methods, it looks at the computed flow and estimates its plausibility given their learned model. Local flow from  $N \times N$  patches is modeled as a multivariate Gaussian; the parameters of this model are partitioned into center flow and the rest. During testing, the center flow vector is evaluated against the rest of the flow in the window, and its correlation based on the trained model is used to determine confidence:

$$\psi_{pVal}^i = \text{inf}\{\alpha \in [0, 1] | d_M(\mathbf{f}_i) > G^{-1}(1 - \alpha)\}, \quad (16)$$

where  $d_M(\mathbf{f}_i)$  is the Mahalanobis distance between the central flow vector  $\mathbf{f}_i$  and its surrounding region given the learned mean and covariance from the training data.  $G^{-1}$  is the inverse of the cumulative distribution function of the distances  $d_M(\cdot)$  obtained from the training data. Our results are provided using our own implementation of their work, where we train the model on 5,000 patches from several sequences (32 in all) with a patch size of 11. As suggested in their paper, we rotate each patch four times by 90 degrees to get a zero mean estimate of the flow.

For each of the competing confidence measurements, we normalize their outputs between 0 and 1.

## 7 EXPERIMENTS

The online Middlebury Optical Flow Evaluation [1] currently ranks over 60 algorithms. We chose the four algorithms which ranked highly on test data and with implementations

available at the time of writing: [51], [52], [3], and [53]. For brevity, we will refer to them as TV, FL, CN, and LD, respectively. The algorithms were used with their default or most successful published settings, though, in practice, the same algorithm with different parameters could be evaluated by our classifier. For quantitative evaluation, we use the average End Point Error (aEPE) metric [54]:

$$aEPE = \frac{1}{N} \sum_i \sqrt{(u_i - u_i^{GT})^2 + (v_i - v_i^{GT})^2}, \quad (17)$$

which equates to the average of all the distances in pixels between each flow vector and the ground truth. Early experimentation with the average angular error [8] produced similar results.

We do leave-one-out evaluation with ground-truth flow from three sources: the eight Middlebury training sequences [1], two Middlebury-like sequences from [39], and 22 of our own synthetic sequences (denoted with an \*, as described in Section 5 with two additions), for a total of 32. In line with the evaluation of [1], the reported scores are the aEPE across the whole image. Error is not reported for areas known to have no flow and for a 10 pixel boundary region around the border of the image. During leave-one-out tests, we also omit any training scenes if they resemble the test scene. Although we evaluate ourselves on data from other sources, we only train on 20 synthetic ground-truth scenes produced by our system.

We have two sections of experiments. In the first section, we evaluate our confidence measure for each of the individual flow algorithms, and compare against other alternative methods. In the second section, we estimate the best combination of optical flow algorithms.

## 7.1 Optical Flow Confidence

We evaluate our algorithm for several different values of error threshold,  $\epsilon_{epe}^s = \{0.1, 0.25, 0.5, 2, 10\}$ , across 32 test sequences. A subset of these results is presented in Fig. 3. We refer the reader to the supplementary material, which can be found in the Computer Society Digital Library at <http://doi.ieeecomputersociety.org/10.1109/TPAMI.2012.171>, for further images. The image displays the confidence results for four different optical flow algorithms for three different sequences: two Middlebury (one real and one synthetic) and one of our own scenes. Each plot displays the aEPE (Y-axis) as a result of removing pixels in order of confidence. So, at 90 percent, we reject the 10 percent we are least confident about and compute the aEPE on the remaining data. The Kway curve shows the confidence as a result of the algorithm selection experiments in Section 7.2, where for each pixel we report the confidence of the winning flow algorithm as determined by our classifier with  $\epsilon_{epe}^s = 2.0$ . For comparison, we also display the optimal ordering, which serves as a lower bound on the best achievable error. We can see from the figure that the confidence measures for different values of  $\epsilon_{epe}^s$  all produce the same downward trend, with the exception of the TV and RubberWhale pairing for  $\epsilon_{epe}^s = 10$ . This can be explained by the fact that the largest magnitude flow vector for this sequence is on the order of 2-3 pixels, i.e., is much lower than the trained value of 10 (the aEPE for the different

algorithms are presented in Table 2). Similarly,  $\epsilon_{epe}^s = 10$  performs best for street1txtr1 due to the large motion in that scene.

### 7.1.1 Comparison to Other Methods

We also compare our results to the other general purpose confidence measures outlined in Section 6. Results for three sequences are presented in Fig. 4 using the same sparsification technique from Fig. 3. Our confidence measure is illustrated at a value of  $\epsilon_{epe}^s = 0.25$ . As can be seen for all three sequences A-C, our confidence measure gives the most consistent performance, always reducing the aEPE as more pixels are removed. We consistently produce better scores when compared to the other measures, with the exception of LD RuberWhale. One explanation for that result is that  $\epsilon_{epe}^s = 0.25$  is not a sensitive enough error threshold for the small errors ( $< 0.1$  pixels) produced by the different algorithms on this sequence. A more appropriate value of  $\epsilon_{epe}^s$  would be 0.1 or less, and as we can see from Fig. 3A, LD this produces a better sparsification curve. It is worth noting that the  $\psi_{pVal}$  measure of [5] produces incorrect results for C street1Txxtr1\* even though it has observed flow patches from this scene in training.

In addition to the qualitative comparison from Fig. 4, we also perform a quantitative comparison to the competing methods. Table 1 contains the aEPE scores for each of the different confidence measures averaged across the 32 test sequences from Table 2 for each flow algorithm. To quantify the success in removing the bad flow vectors, we remove pixels based on the confidence and compute the aEPE for the remaining pixels, averaging across all the sequences. For each confidence measure, we evaluate the aEPE at  $P_{amt} = \{30, 60, 90\}$  pixels, where  $P_{amt}$  is the percentage of remaining pixels. As can be observed in Fig. 4, there are instances where there are no pixels remaining at a particular value of  $P_{amt}$ . This is because, in certain situations, multiple pixels can have the exact same confidence value. If there are no pixels within 10 percent of the desired value of  $P_{amt}$ , we simply ignore that aEPE when computing the total average. Our confidence measure produces the best overall results of all the competing methods.

### 7.1.2 Confidence in X and Y Directions

For view interpolation or panoramic stitching in the presence of moving objects, one component of the flow vector could be quite accurate while the other is highly uncertain. In these applications, it could be useful to distinguish this situation from one where both X and Y flow is unconfident. In addition to computing a joint confidence, we can also produce a confidence for the horizontal and vertical directions separately. We simply train the same classifier on either the X or Y flow components. Fig. 5 shows the separated confidence images for two scenes, one featuring horizontal motion and the other vertical. In the first sequence, we can see that our confidence measure is more confident for flow vectors in the Y direction (as there is very little to no vertical motion) but more uncertain in the X direction. The second scene depicts several objects falling to the ground. Our Y confidence correctly identifies more uncertainty in the vertical direction.

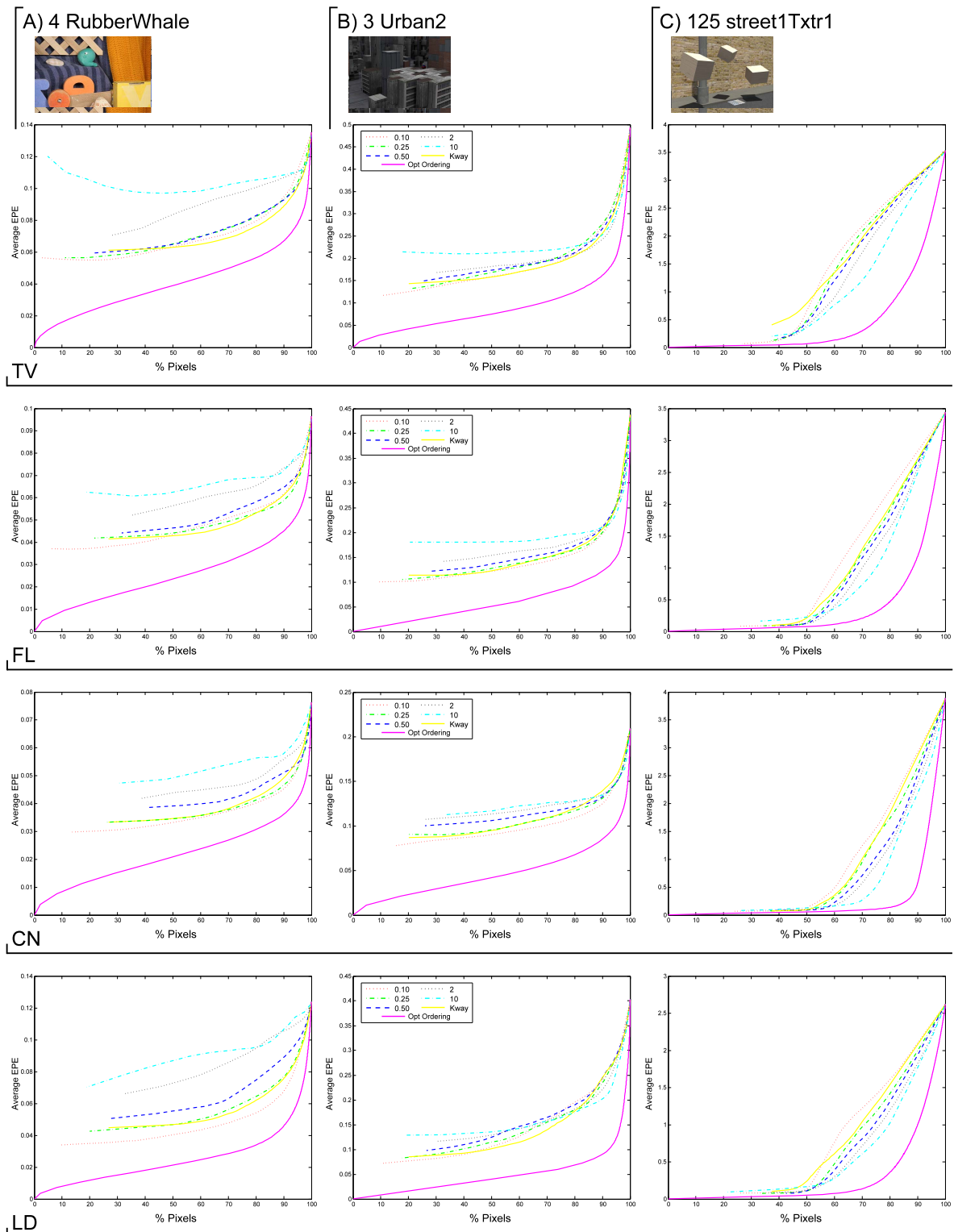


Fig. 3. Confidence graphs. Each row represents a different algorithm, while each column is one of three different scenes. Our confidence measure is illustrated at different values of the error threshold,  $\epsilon_{epe}^s = \{0.1, 0.25, 0.5, 2, 10\}$ . Kway represents the confidence for the combined flow using  $\epsilon_{epe}^s = 2.0$ . Each scene/algorithm pair displays the aEPE as a result of keeping  $x$  percent of flow vectors in order of diminishing confidence. Note that the Y-axis for each scene/algorithm pair has a different scaling.

## 7.2 Choosing the Best Optical Flow Algorithm

In our next set of experiments, we predict which one of  $K$  constituent optical flow algorithms (in this case,  $K = 4$ ; TV, FL, CN, LD) to trust at each pixel. We perform leave-one-out tests on all 32 sequences. The results are summarized in

Table 2. It is interesting to note that of the four algorithms, none outperforms the others on all sequences. While CN gives the best results on more sequences, LD has a lower aEPE overall. “OptCombo” is the optimal combination given the ground truth, and serves as a lower limit on the

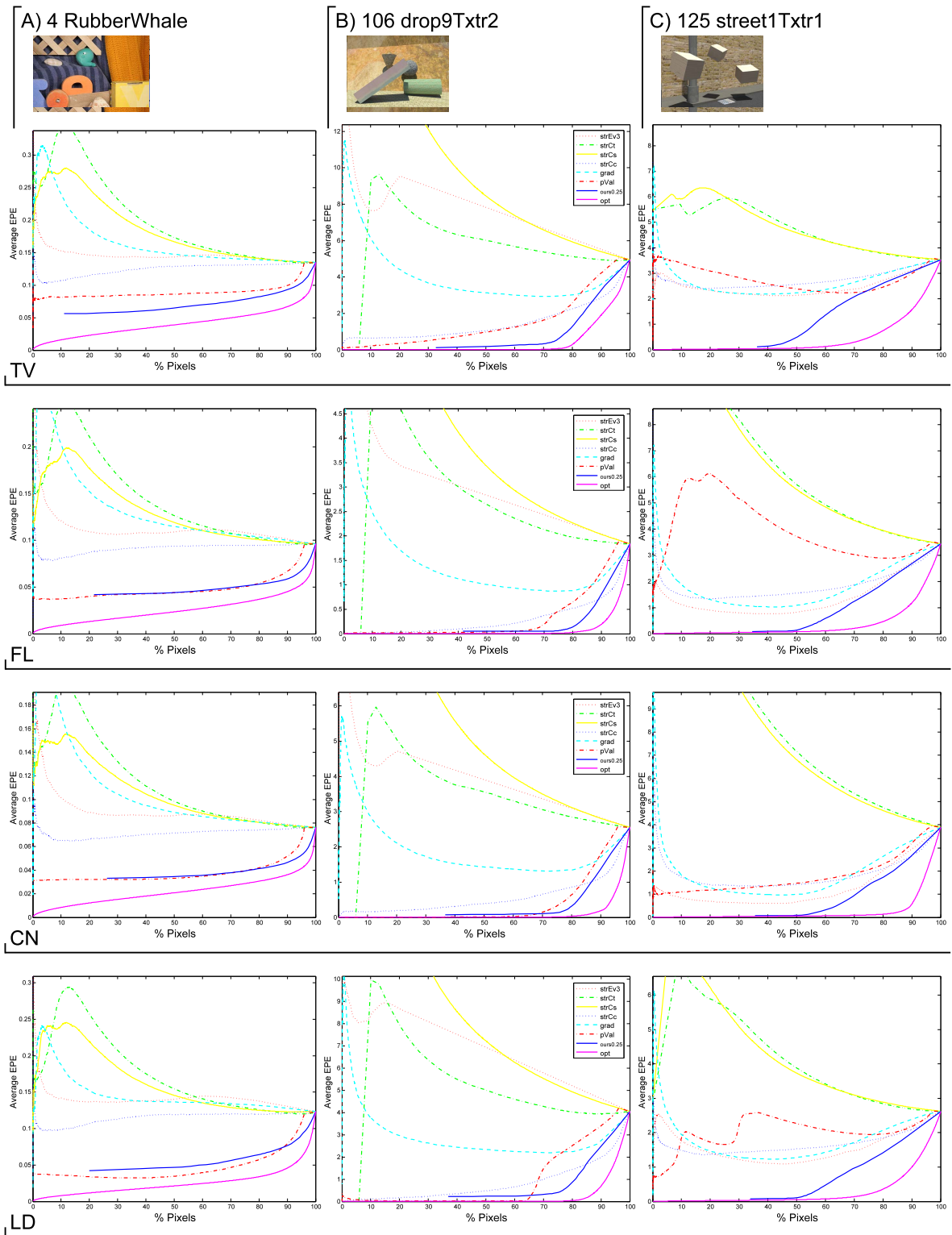


Fig. 4. Confidence comparison. We compare our confidence measure against six others. This image presents three different scenes and four different optical flow algorithms. We follow the same sparsification technique as shown in Fig. 3. Our measure, “ours0.25,” consistently ranks the flow vectors by accuracy better than any other method.

best possible performance achievable. “OursKWay” is the multiclass formulation from (3). “OursCombo” combines the output of the  $K$  individual confidence measures for each flow algorithm. At each pixel, we choose the algorithm that is the most confident. The results here are presented for  $\epsilon_{epe}^s = 2.0$ . From Fig. 3, we can see that almost any other value of  $\epsilon_{epe}^s$

would also perform very well. Interestingly, it does not win any of the individual sequences, but it consistently comes a close second and achieves the best score overall. “OursScene” chooses the result of the algorithm which has the majority vote from the classifier. “RandCombo” is a baseline algorithm that simply randomly chooses one of the  $K$  algorithms at

TABLE 1  
Confidence Measure Comparison

method	TV	FL	CN	LD
strEv3	1.525	1.548	1.347	1.070
strCt	1.824	1.978	1.813	1.184
strCs	1.708	1.739	1.579	1.229
strCc	1.400	1.323	1.283	0.912
grad	1.371	1.622	1.423	0.887
pVal	1.137	0.780	0.829	0.831
ours $\epsilon_{epe}1.0$	0.605	0.504	<b>0.568</b>	<b>0.416</b>
ours $\epsilon_{epe}0.25$	<b>0.512</b>	<b>0.381</b>	0.600	0.453

Each of the competing confidence measures is evaluated on the different flow algorithms across the 32 test sequences. Each score represents the total average computed by removing a different number of pixels and only counting the score for the remaining  $P_{amt} = \{30, 60, 90\}$ , with lower scores being better.

each location; as expected, it performs the worst overall. It is worth noting that the best result could bear improving to further close in on the ideal possible combination.

In sequence 24 (Crates2Httrx1\*), the different flow algorithms produce widely varying aEPE scores (see Table 2). In Fig. 6, we can see the results of our algorithm selection for “OursKWay.” Our classifier avoids regions of large error, which is most noticeable on the blue crate in the foreground. Instead of choosing the flow estimated by CN (which generally gives very good performance), it chooses flow from LD and FL. Here, the color coding can be slightly misleading as it simply shows the flow algorithm with the

highest probability and does not indicate how close the different algorithms actually are. Both our methods for predicting the best combination produce results close to the winning result for this particular scene.

### 7.2.1 Effects of Training Data

Fig. 7 illustrates the effect of varying the amount of training data used from each sequence while doing leave-one-out tests on two sequences. To minimize the effects of randomness, we ran each of these leave-one-out experiments three times and averaged their results. For both scenes, we can see that the aEPE slightly improves as more data are included in training. This is explained by the fact that for most sequences it is very difficult to reduce the aEPE. Typically, regions such as object boundaries contain most of the error. A more revealing metric is to look at the aEPE of the flow vectors with the worst accuracy (in this example, we look at the worst 5 percent of vectors—aveEpeTp5). Sequence 4 (RubberWhale) has only a slight improvement. This is because each of the constituent flow algorithms has a similar aEPE, see Table 2, whereas, 88 (blow19Ttrx2\*) benefits from more training data.

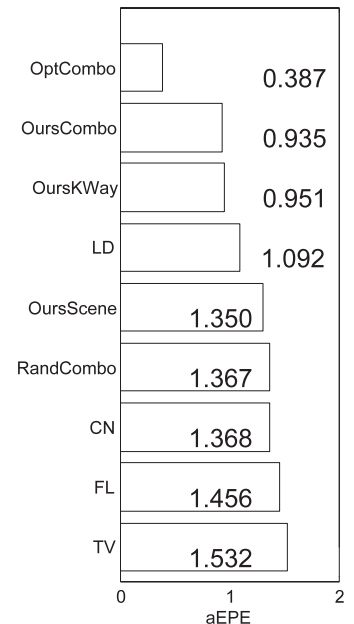
Fig. 8 illustrates the feature importance as given by the Random Forest classifier for sequence 24 (Crates2Httrx1\*) for the “OursKWay” leave-one-out experiment from Table 2. We see the boundary features are the most important, followed by the temporal gradients.

TABLE 2

Leave-One-Out Average EPE Scores for Estimating the Best Optical Flow Algorithm at Each Pixel Location for 32 Different Scenes

Image Sequence	TV	FL	CN	LD	OursKWay	OursCombo	RandCombo	OptCombo
1 Venus	0.408	0.342	<b>0.229</b>	0.433	0.304	0.306	0.351	0.176
2 Urban3	1.132	0.524	<b>0.377</b>	0.600	0.502	0.543	0.658	0.200
3 Urban2	0.506	0.444	<b>0.207</b>	0.334	0.353	0.331	0.368	0.123
4 RubberWhale	0.135	0.096	<b>0.077</b>	0.120	0.092	0.108	0.107	0.052
5 Hydrangea	0.196	0.164	<b>0.154</b>	0.178	0.169	0.168	0.174	0.100
6 Grove3	0.745	0.624	<b>0.438</b>	0.657	0.605	0.585	0.616	0.324
7 Grove2	0.220	0.169	<b>0.091</b>	0.159	0.149	0.161	0.159	0.064
8 Dimetrodon	0.211	0.144	<b>0.115</b>	0.117	0.139	0.152	0.147	0.077
9 Crates1*	3.464	3.730	3.150	<b>3.104</b>	3.234	3.113	3.365	2.423
10 Crates2*	4.615	12.572	10.409	<b>2.513</b>	2.617	3.692	7.568	1.544
13 Mayan1*	2.331	<b>0.727</b>	1.718	5.567	3.887	2.590	2.626	0.297
14 Mayan2*	0.442	0.344	<b>0.211</b>	0.350	0.304	0.247	0.339	0.138
15 YosemiteSun	0.310	0.250	0.232	<b>0.188</b>	0.251	0.253	0.245	0.142
16 GroveSun	0.576	0.403	<b>0.233</b>	0.484	0.301	0.335	0.424	0.170
17 Robot*	2.335	1.857	1.525	1.212	<b>1.005</b>	1.133	1.734	0.415
18 Sponzal*	1.006	1.013	1.102	<b>0.917</b>	1.009	0.997	1.010	0.635
19 Sponza2*	0.531	0.494	1.674	<b>0.481</b>	1.538	1.485	0.791	0.307
22 Crates1Httrx2*	1.106	0.693	1.640	<b>0.548</b>	0.931	0.679	0.999	0.210
24 Crates2Httrx1*	3.128	10.210	8.805	<b>0.809</b>	1.222	2.080	5.762	0.382
26 Brickbox1t1*	1.094	0.394	<b>0.228</b>	2.602	0.373	0.457	1.070	0.148
29 Brickbox2t2*	7.478	1.827	2.192	3.505	<b>1.690</b>	1.802	3.765	0.716
30 GrassSky*	2.102	2.484	1.317	<b>1.039</b>	1.750	1.209	1.746	0.434
39 GrassSky9*	0.722	0.438	<b>0.273</b>	0.510	0.378	0.358	0.486	0.189
49 TxtRMovement*	3.166	0.241	<b>0.132</b>	0.356	0.337	0.331	0.969	0.063
50 TxtLMovement*	1.521	0.282	<b>0.126</b>	0.604	0.225	0.318	0.652	0.057
51 blow1Ttrx1*	0.085	0.050	<b>0.027</b>	0.081	0.048	0.052	0.061	0.017
88 blow19Ttrx2*	0.525	0.380	<b>0.199</b>	0.319	0.301	0.311	0.355	0.145
89 drop1Ttrx1*	0.119	0.071	<b>0.052</b>	0.084	0.063	0.070	0.082	0.026
106 drop9Ttrx2*	5.195	<b>1.985</b>	2.715	4.369	3.301	3.095	3.574	1.362
107 roll1Ttrx1*	0.004	0.005	<b>0.002</b>	0.002	0.003	0.004	0.003	0.001
124 roll9Ttrx2*	0.040	0.048	<b>0.014</b>	0.023	0.022	0.027	0.031	0.011
125 street1Ttrx1*	3.647	3.585	4.097	<b>2.664</b>	3.329	2.923	3.494	1.446

TV, FL, CN, and LD are the four constituent optical flow algorithms. “OursKWay” is the result of our multiclass classification, “OursCombo” is the combined best confidence, “RandCombo” is a random combination, “OursScene” is the algorithm with the most votes for a scene, and “OptCombo” represents the optimal ground-truth combination. The bar chart to the right displays total aEPE across all 32 sequences.



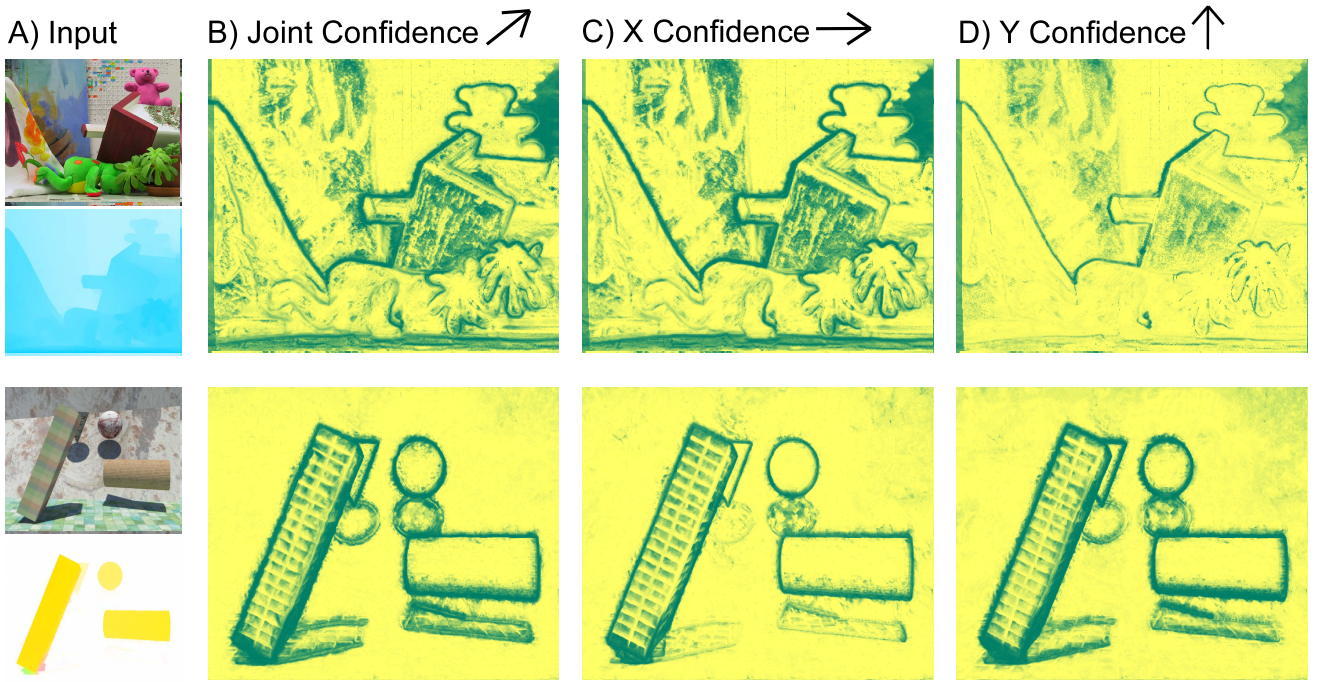


Fig. 5. Horizontal and vertical confidence. A) Input image and computed flow. B) Estimated confidence. C) Estimated confidence in  $X$  direction. D) Estimated confidence in  $Y$  direction. Each row represents a different scene with confidence computed for CN [3]. The first scene features predominantly horizontal motion and is from the Middlebury Stereo dataset with  $\epsilon^s_{epe} = 0.3$ . It can be seen that the  $Y$  confidence image is more certain than the joint confidence. The second scene, 89 drop1Txxr1\*, features vertical motion with  $\epsilon^s_{epe} = 0.1$ . There is more uncertainty around the falling objects in the  $Y$  confidence as compared to  $X$ .

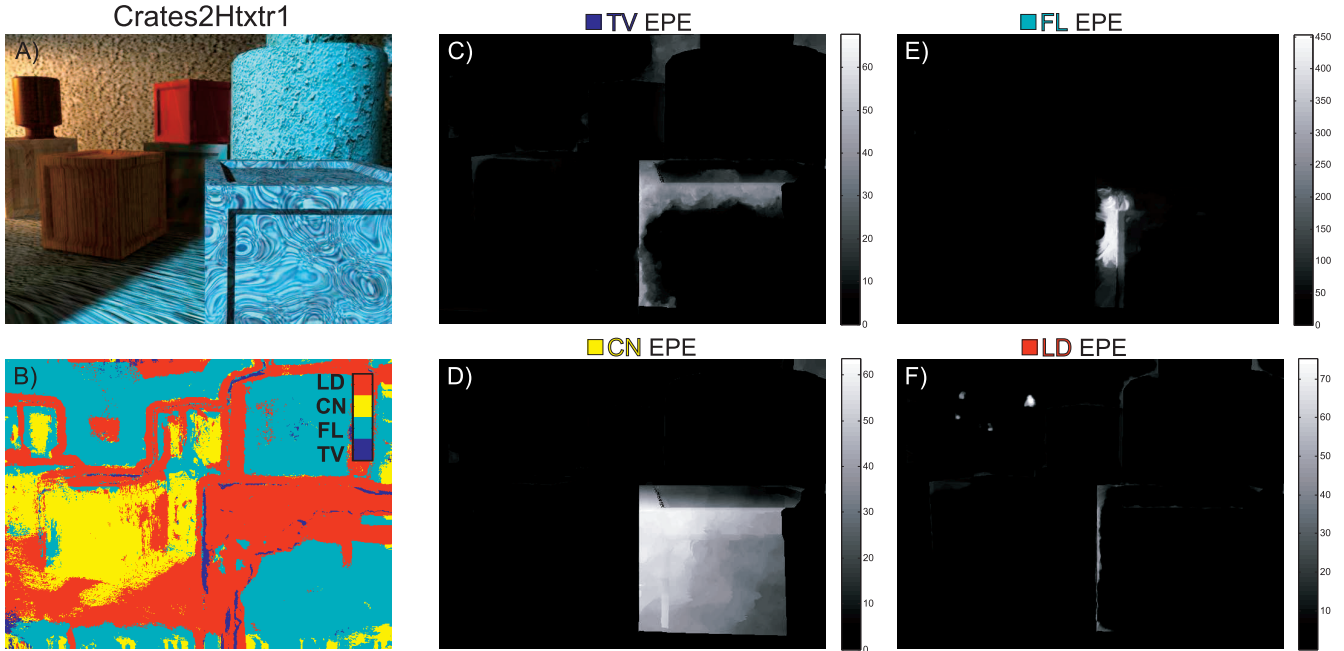


Fig. 6. Selecting the best algorithm. A) First image of input pair. B) Selected algorithm result for “OursKWay” classification where each color represents a different algorithm. C), D), E), F) EPE images for TV, CN, FL, and LD. Note how, in different image regions, our classifier avoids choosing flow algorithms which produce larger errors.

### 7.3 Implementation Details

For all experiments, the Random Forest classifier was run with 50 trees, 50 minimum samples at each node, and a maximum tree depth of 10. We use all our scenes from Table 2 marked with an \* for training, with the exception of 13 and 14, which are omitted due to their low resolution. In total, we have 20, 640 × 480, training sequences with 14,000 samples randomly chosen from each and with an even

amount for each class in the suitability experiments. For the feature vector, the hysteresis threshold  $\tau_{ed}$  is set to the value returned by Matlab, and for  $\tau_{pb}$  we use 0.1 and 0.4. For the photoconstancy residual, we set a value of 1,000 if the flow vector points out of the frame. For the features that exploit scale, we use an image pyramid with  $z = [1, 10]$  levels and a rescaling factor of  $s = 0.8$ . Due to their computational expense, the  $P_b$  features were computed for four levels.

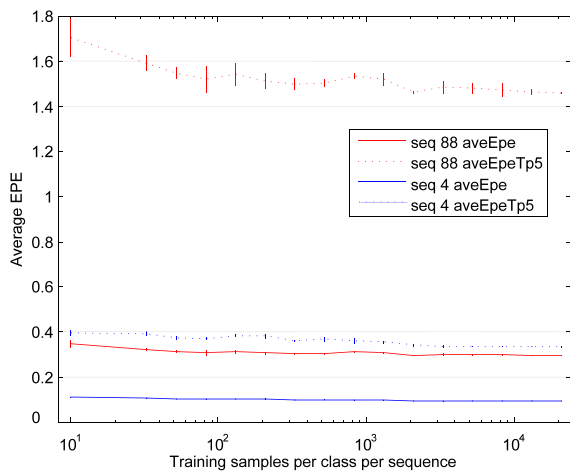


Fig. 7. The effect on average EPE for “OursKWay” on two different sequences as we change the number of training samples per class per sequence. aveEpe is the average EPE and aveEpeTp5 is the average EPE for the worst 5 percent of the data. Error bars show the standard deviation.

Combining all the features results in a 52-dimensional feature vector.

Our unoptimized code for the classifier is implemented in C and features are computed in Matlab. All the following times are presented in seconds for a typical  $640 \times 480$  image pair on an Intel i7 2.67 GHz with 6 GB RAM. Computation for each of the flow algorithms takes in seconds: TV: 35.5, FL: 22.4, CN: 1,330.8, and LD: 258.2. The Random Forest takes 439.8 seconds to train on 20 sequences and testing takes 30 seconds. Our features are all quite lightweight and take 6.5 seconds to compute, not including the  $P_b$  features, which take a total of 2,165 seconds for four scales. These features could be sped up using a more efficient C implementation. The current major bottleneck is the need to compute the different optical flow vectors. With more GPU implementations for flow (e.g., [52]) these times will hopefully reduce.

## 8 CONCLUSIONS

There is an ever-increasing variety of solutions to the problem of optical flow estimation. These different algorithms and their energy functions can be seen as good or bad, but only with respect to specific video situations. Our main finding is that the success (or failure) of all the flow algorithms we tested for aEPE is predictable, given our supervised learning framework.

Each algorithm processes sequences differently. Frames encoded with our feature vector (9) correlate well with the applicability of that process for each sequence. Our feature vector embodies two important characteristics. First, it is comprised of multiple different measures, incorporating a broad range of motion and appearance cues and simple algorithm-specific qualities like the photoconstancy residual. Second, mapping feature vectors to uncertainty labels using a Random Forest means that the training process performs feature selection. Instead of heuristic choices about the expected smoothness of flow fields or anticipated challenges of textureless regions, our method objectively chooses weights, picking out which features are important and in what combinations.

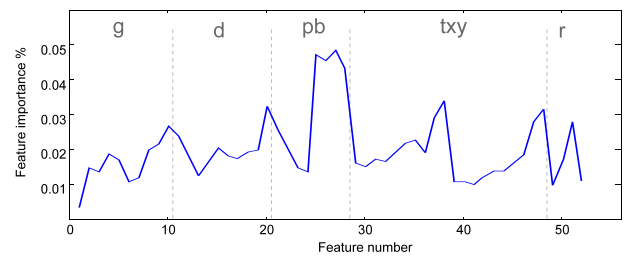


Fig. 8. Feature importance as given by the random forest classifier for sequence 24 (Crates2Htxtr1\*) for the “OursKWay” leave-one-out experiment from Table 2.

Per-algorithm flow confidence is worth measuring, and can be applied to whole videos or just parts. The OptCombo and OursCombo columns of Table 2 show that even though modern algorithms agree on much of a scene’s flow, significant disagreements are worth settling by carefully modeling each algorithm’s uncertainty. Knowing where a flow algorithm’s performance is predicted to be uncertain creates opportunities for interesting applications. We have shown (Fig. 3) that excluding pixels for which the flow confidence is low really reduces the overall aEPE. The impact is different on different sequences, sometimes by an order of magnitude, but consistently improves performance for  $aEPE < 2$  pixels. Now a user of an existing or future flow algorithm can balance their need for spatial coverage (i.e., number of pixels) against the accuracy they can accept. Further, they can decide to keep or ignore flow which is only confident in the X or Y direction, allowing for higher level algorithms to degrade gracefully when full 2D flow is underconstrained. Finally, whether using our multiclass or one-versus-all methods, users can now elect to locally (see Fig. 6) trust each algorithm only where it is most appropriate. We advocate this resulting algorithm-suitability “patchwork” for users who want the lowest EPE over videos in general, and who, for specific videos and with little guesswork, prefer to consistently win “silver” instead of an occasional “gold.”

### 8.1 Limitations and Future Work

The most exciting avenue for future work is the opportunity to develop specialist flow algorithms for narrowly defined situations. If the situation can be detected using our framework, then that specialist algorithm could be terrible in general, as long as it excels in its narrow domain. Opportunities obviously exist for further features that may also correlate with flow confidence. Heterogeneous features that incorporate context of motion could be quite revealing, and more accurate occlusion information [19] could prove useful.

One limitation of Random Forests and most supervised learning algorithms is that a training example specifies only that one algorithm is most suitable, while the rest are equally unsuitable. This effectively ignores the fact that the second-best algorithm could give an end-point estimate 10 times closer than the fourth best. Equally, when differences between the top two algorithms are minimal, we must currently either ignore the example completely or expend effort trying to learn to distinguish between equals. Our one-versus-all tests were posed as classification

challenges to allow easy comparisons between experiments, but a system similar to our prototype could be built around regressing per-algorithm flow-confidence. Using regression instead of classification would potentially allow us to automatically choose the best value of error threshold  $\epsilon_{epe}^s$  instead of relying on a user-provided value. This could allow us to learn the relationship between image features and optical flow error directly.

We chose to validate our approach on the example problem of optical flow. There are many other applications, such as stereo, where multiple competing algorithms vie to be universally best, and it would be interesting to try our learned segmentation approach there. Also, a reliable estimate of confidence for these problems would be of great use to practitioners [55]. Finally, our approach ignores the cost of processing times, which is currently acceptable, but  $O(k)$  in the number of algorithms under consideration. One strategy could be to optimize the classifier subject to the computational cost of each algorithm.

## ACKNOWLEDGMENTS

Funding for this research was provided by the National University of Ireland Travelling Studentship in the Sciences and the Microsoft Innovation Cluster for Embedded Software. Project page: <http://visual.cs.ucl.ac.uk/pubs/flowConfidence>.

## REFERENCES

- [1] S. Baker, D. Scharstein, J.P. Lewis, S. Roth, M.J. Black, and R. Szeliski, "A Database and Evaluation Methodology for Optical Flow," *Int'l J. Computer Vision*, vol. 92, pp. 1-31, 2011.
- [2] O. Mac Aodha, G.J. Brostow, and M. Pollefeys, "Segmenting Video into Classes of Algorithm-Suitability," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2010.
- [3] D. Sun, S. Roth, and M. Black, "Secrets of Optical Flow Estimation and Their Principles," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2010.
- [4] J. Kybic and C. Nieuwenhuis, "Bootstrap Optical Flow Confidence and Uncertainty Measure," *Computer Vision and Image Understanding*, vol. 115, no. 10, pp. 1449-1462, 2011.
- [5] C. Kondermann, R. Mester, and C. Garbe, "A Statistical Confidence Measure for Optical Flows," *Proc. European Conf. Computer Vision*, 2008.
- [6] S. Gould, R. Fulton, and D. Koller, "Decomposing a Scene into Geometric and Semantically Consistent Regions," *Proc. IEEE Int'l Conf. Computer Vision*, 2009.
- [7] A. Farhadi, I. Endres, D. Hoiem, and D. Forsyth, "Describing Objects by Their Attributes," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2009.
- [8] J.L. Barron, D.J. Fleet, and S.S. Beauchemin, "Performance of Optical Flow Techniques," *Int'l J. Computer Vision*, vol. 12, pp. 43-77, 1994.
- [9] E. Simoncelli, E. Adelson, and D. Heeger, "Probability Distributions of Optical Flow," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 1991.
- [10] P. Anandan, "A Computational Framework and an Algorithm for the Measurement of Visual Motion," *Int'l J. Computer Vision*, vol. 2, pp. 283-310, 1989.
- [11] S. Uras, F. Girosi, A. Verri, and V. Torre, "A Computational Approach to Motion Perception," *Biological Cybernetics*, vol. 60, pp. 79-87, 1988.
- [12] B. Jähne, H. Haussecker, and P. Geissler, *Handbook of Computer Vision and Applications: Signal Processing and Pattern Recognition*, vol. 2. Academic Press, 1999.
- [13] A. Wedel, D. Cremers, T. Pock, and H. Bischof, "Structure- and Motion-Adaptive Regularization for High Accuracy Optic Flow," *Proc. IEEE Int'l Conf. Computer Vision*, 2009.
- [14] H. Zimmer, A. Bruhn, J. Weickert, L. Valgaerts, A. Salgado, B. Rosenhahn, and H.-P. Seidel, "Complementary Optic Flow," *Proc. Seventh Int'l Conf. Energy Minimization Methods in Computer Vision and Pattern Recognition*, 2009.
- [15] A. Bruhn and J. Weickert, "A Confidence Measure for Variational Optic Flow Methods," *Geometric Properties for Incomplete Data*, pp. 283-298, Springer, 2006.
- [16] J.L. Barron, D.J. Fleet, and S.S. Beauchemin, "Performance of Optical Flow Techniques," Technical Report TR299, Dept. of Computer Science, Univ. of Western Ontario, 1992.
- [17] M. Gong and Y.-H. Yang, "Estimate Large Motions Using the Reliability-Based Motion Estimation Algorithm," *Int'l J. Computer Vision*, vol. 68, pp. 319-330, 2006.
- [18] C. Kondermann, D. Kondermann, B. Jähne, and C. Garbe, "An Adaptive Confidence Measure for Optical Flows Based on Linear Subspace Projections," *Proc. DAGM Conf. Pattern Recognition*, pp. 132-141, 2007.
- [19] A. Humayun, O. Mac Aodha, and G.J. Brostow, "Learning to Find Occlusion Regions," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2011.
- [20] S. Gehrig and T. Scharwächter, "A Real-Time Multi-Cue Framework for Determining Optical Flow Confidence," *Proc. IEEE Int'l Conf. Computer Vision Workshops*, 2011.
- [21] A. Bainbridge-Smith and R. Lane, "Measuring Confidence in Optical Flow Estimation," *Electronics Letters*, vol. 32, no. 10, pp. 882-884, 1996.
- [22] M. Reynolds, J. Dobos, L. Peel, T. Weyrich, and G. Brostow, "Capturing Time-of-Flight Data with Confidence," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2011.
- [23] B. Li, R. Xiao, Z. Li, R. Cai, B.-L. Lu, and L. Zhang, "Rank-Sift: Learning to Rank Local Interest Points," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2011.
- [24] V.C. Raykar, S. Yu, L.H. Zhao, A. Jerebko, C. Florin, G.H. Valadez, L. Bogoni, and L. Moy, "Supervised Learning from Multiple Experts: Whom to Trust When Everyone Lies a Bit," *Proc. Int'l Conf. Machine Learning*, 2009.
- [25] X. Yong, D. Feng, Z. Rongchun, and M. Petrou, "Learning-Based Algorithm Selection for Image Segmentation," *Pattern Recognition Letters*, vol. 26, pp. 1059-1068, 2005.
- [26] B. Stenger, T. Woodley, and R. Cipolla, "Learning to Track with Multiple Observers," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2009.
- [27] N. Alt, S. Hinterstoisser, and N. Navab, "Rapid Selection of Reliable Templates for Visual Tracking," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2010.
- [28] B. Peng and O. Veksler, "Parameter Selection for Graph Cut Based Image Segmentation," *Proc. British Machine Vision Conf.*, 2008.
- [29] M. Muja and D.G. Lowe, "Fast Approximate Nearest Neighbors with Automatic Algorithm Configuration," *Proc. Int'l Conf. Computer Vision Theory and Applications*, 2009.
- [30] D.G. Lowe, "Distinctive Image Features from Scale-Invariant Keypoints," *Int'l J. Computer Vision*, vol. 60, pp. 91-110, 2004.
- [31] V. Lempitsky, S. Roth, and C. Rother, "Fusionflow: Discrete-Continuous Optimization for Optical Flow Estimation," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2008.
- [32] B.D. Lucas and T. Kanade, "An Iterative Image Registration Technique with an Application to Stereo Vision," *Proc. Int'l Joint Conf. Artificial Intelligence*, 1981.
- [33] B. Horn and B.G. Schunck, "Determining Optical Flow," *Artificial Intelligence*, vol. 17, pp. 185-204, 1981.
- [34] A. Bruhn, J. Weickert, and C. Schnörr, "Lucas/Kanade Meets Horn/Schunck: Combining Local and Global Optic Flow Methods," *Int'l J. Computer Vision*, vol. 61, pp. 211-231, 2005.
- [35] S.C. Zhu, Y.N. Wu, and D. Mumford, "Filters, Random Fields and Maximum Entropy (FRAME): Towards a Unified Theory for Texture Modeling," *Int'l J. Computer Vision*, vol. 27, pp. 107-126, 1998.
- [36] S. Roth and M.J. Black, "Fields of Experts," *Int'l J. Computer Vision*, vol. 82, pp. 205-229, 2009.
- [37] O. Woodford, I.D. Reid, P.H.S. Torr, and A.W. Fitzgibbon, "Fields of Experts for Image-Based Rendering," *Proc. British Machine Vision Conf.*, 2006.
- [38] S. Roth and M.J. Black, "On the Spatial Statistics of Optical Flow," *Int'l J. Computer Vision*, vol. 74, pp. 33-50, 2007.
- [39] D.Q. Sun, S. Roth, J.P. Lewis, and M.J. Black, "Learning Optical Flow," *Proc. European Conf. Computer Vision*, 2008.

- [40] A. Levin, D. Lischinski, and Y. Weiss, "Colorization Using Optimization," *Proc. ACM Siggraph*, 2004.
- [41] L. Breiman, "Random Forests," *Machine Learning*, vol. 45, pp. 5-32, 2001.
- [42] A. Criminisi, J. Shotton, and E. Konukoglu, "Decision Forests: A Unified Framework for Classification, Regression, Density Estimation, Manifold Learning and Semi-Supervised Learning," *Foundations and Trends in Computer Graphics and Vision*, vol. 7, pp. 81-227, 2012.
- [43] D. Martin, C. Fowlkes, and J. Malik, "Learning to Detect Natural Image Boundaries Using Local Brightness, Color, and Texture Cues," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 26, no. 5, pp. 530-549, May 2004.
- [44] C. Liu, W.T. Freeman, E.H. Adelson, and Y. Weiss, "Human-Assisted Motion Annotation," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2008.
- [45] J. Shotton, A. Fitzgibbon, M. Cook, T. Sharp, M. Finocchio, R. Moore, A. Kipman, and A. Blake, "Real-Time Human Pose Recognition in Parts from Single Depth Images," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2011.
- [46] O. Mac Aodha, N.D.F. Campbell, A. Nair, and G.J. Brostow, "Patch Based Synthesis for Single Depth Image Super-Resolution," *Proc. European Conf. Computer Vision*, 2012.
- [47] J. Cech, J. Sanchez-Riera, and R. Horaud, "Scene Flow Estimation by Growing Correspondence Seeds," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2011.
- [48] B. McCane, K. Novins, D. Crannitch, and B. Galvin, "On Benchmarking Optical Flow," *Computer Vision and Image Understanding*, vol. 84, pp. 126-143, 2001.
- [49] B. Kaneva, A. Torralba, and W. Freeman, "Evaluation of Image Features Using a Photorealistic Virtual World," *Proc. IEEE Int'l Conf. Computer Vision*, 2011.
- [50] S. Meister and D. Kondermann, "Real Versus Realistically Rendered Scenes for Optical Flow Evaluation," *Proc. 14th ITG Conf. Electronic Media Technology*, 2011.
- [51] C. Zach, T. Pock, and H. Bischof, "A Duality Based Approach for Realtime TV-L1 Optical Flow," *Proc. DAGM Conf. Pattern Recognition*, 2007.
- [52] M. Werlberger, W. Trobin, T. Pock, A. Wedel, D. Cremers, and H. Bischof, "Anisotropic Huber-L1 Optical Flow," *Proc. British Machine Vision Conf.*, 2009.
- [53] T. Brox and J. Malik, "Large Displacement Optical Flow: Descriptor Matching in Variational Motion Estimation," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 33, no. 3, pp. 500-513, Mar. 2011.
- [54] M. Otte and H. Nagel, "Optical Flow Estimation: Advances and Comparisons," *Proc. European Conf. Computer Vision*, 1994.
- [55] X. Hu and P. Mordohai, "A Quantitative Evaluation of Confidence Measures for Stereo Vision," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 34, no. 11, pp. 2121-2133, Nov. 2012.



**Oisín Mac Aodha** received the BEng degree in electronic and computer engineering from the National University of Ireland in Galway and the MSc degree in Machine Learning from University College London (UCL). In 2009, he worked as a research assistant in Marc Pollefeys group at ETH Zurich. In 2010, he was awarded an NUI Travelling Studentship in the Sciences and is currently working toward the PhD degree in the Computer Science Department at UCL. His research interests lie in the application of machine learning to 3D scene processing. He is a student member of the IEEE.

learning to video data, and discrete optimization techniques for video segmentation. He is a student member of the IEEE.



**Ahmad Humayun** received the MS degree in Computer Graphics, Vision, and Imaging from University College London in 2010, for which he won the Microsoft Research Project Prize and BBC Best Overall Student Prize. He joined the Computational Perception Lab at the Georgia Institute of Technology in 2011, where he is currently pursuing the PhD degree in computer vision. His research interests are in the use of motion cues for perception, applying machine

learning to video data, and discrete optimization techniques for video segmentation. He is a student member of the IEEE.



**Marc Pollefeys** received the PhD degree in 1999 from the KULeuven. He has been a full professor in the Department of Computer Science of ETH Zurich since 2007. He also remains associated with the Department of Computer Science of the University of North Carolina at Chapel Hill, where he started as an assistant professor in 2002. His main area of research is geometric computer vision and his aim is to develop flexible approaches to capture visual representations of real-world objects, scenes, and events. He received the Marr prize, a US National Science Foundation (NSF) CAREER award, a Packard Fellowship, and an ERC grant. He is the author of more than 170 peer-reviewed publications. He is the general chair for ECCV '14 and was a program cochair for CVPR '09. He is on the editorial board of the *International Journal of Computer Vision*, is an associate editor for the *IEEE Transaction Pattern Analysis And Machine Intelligence*, and is a fellow of the IEEE.



**Gabriel J. Brostow** received the BS degree in electrical engineering from the University of Texas at Austin in 1996 and the MS and PhD degrees in computer science from the Georgia Institute of Technology in 2004. He was a Marshall Sherfield fellow and postdoctoral researcher at the University of Cambridge until 2007, and a research scientist at ETH Zurich until 2009. In 2008, he joined the Computer Science Department at University College London as an

assistant professor. His research focus is "smart capture" of visual data. He is a member of the IEEE.

► For more information on this or any other computing topic, please visit our Digital Library at [www.computer.org/publications/dlib](http://www.computer.org/publications/dlib).