

Spatiotemporal Stereo via Spatiotemporal Quadric Element (Stequel) Matching

Mikhail Sizintsev and Richard P. Wildes
Department of Computer Science and Engineering
York University
Toronto, ON, Canada

Abstract

Spatiotemporal stereo is concerned with the recovery of the 3D structure of a dynamic scene from a temporal sequence of multiview images. This paper presents a novel method for computing temporally coherent disparity maps from a sequence of binocular images through an integrated consideration of image spacetime structure and without explicit recovery of motion. The approach is based on matching spatiotemporal quadric elements (stequels) between views, as it is shown that this matching primitive provides a natural way to encapsulate both local spatial and temporal structure for disparity estimation. Empirical evaluation with laboratory-based imagery with ground truth and more typical natural imagery shows that the approach provides considerable benefit in comparison to alternative methods for enforcing temporal coherence in disparity estimation.

1. Introduction

In a 3D dynamic environment a visual system must process image data that derives from both the temporal and spatial scene dimensions. Correspondingly, stereo and motion are two of the most widely researched areas in computer vision. Within this body of research, integrated investigation of stereo and motion has received relatively little attention. Ultimately, however, recovery of 3D scene structure must respect dynamic information to ensure that estimates are temporally consistent. Further, in situations where instantaneous multiview matching is ambiguous (e.g., weakly textured surfaces or epipolar aligned pattern structure), dynamic information has the potential to resolve correspondence by further constraining possible matches.

In response to the above observations, this paper describes a novel approach to recovering temporally coherent disparity estimates from a sequence of binocular images. The key idea is to base stereo correspondence on matching primitives that inherently encompass both the spatial and temporal dimensions of image spacetime. In particular, each temporal stream of imagery is locally represented in terms of its orientation structure, as captured by the spatiotemporal quadric (also variously referred to as the orientation tensor and covariance matrix, see, e.g., [10]). It will

be shown that by basing matching on this representation, it is possible to recover temporally coherent disparity estimates, without the need to make optical or 3D flow explicit. Further, this representation allows spatial and temporal image structure to resolve otherwise ambiguous matches in a fashion consistent with both sources of information.

Early work combining stereo and motion concentrated on punctate features (e.g., edges, corners). One of the earliest attempts made use of heuristics for assigning spatial and temporal matches based on model-based reasoning [12]. A rather different early approach exploited constraints on the temporal derivative of disparity [26]. Other work matched binocular features to recover 3D estimates for temporal tracking [29]. More recent research that relies on loose coupling of stereo and motion has emphasized the recovery of 3D motion using optical flow in conjunction with multiple hypothesis disparity maps [5]. The proposed research differs from such early work in being focused on a more integrated approach to spatiotemporal processing and in its emphasis on dense reconstruction.

More recent stereo research has seen increased interest in scene recovery from multicamera (especially binocular) video as constrained by 3D models. Some work has concentrated on the recovery of surface mesh models between individual stereo pairs with tracking across time instances serving to yield temporally consistent models [16]. Other research considers multiple cameras, employs voxel carving for initial estimation and uses intensity-based matching over spatiotemporal volumes without consideration of image motion differences between different views [18]. Still other work casts stereo and motion estimation as a generic image matching problem solved variationally after backprojecting the input images onto a suitable surface [20]. Again, the present work puts emphasis on a more integrated approach to stereo and motion and in eschewing explicit surface models, which can become problematic when dealing with multiple objects and complex scenes.

Other lines of recent research have emphasized more integrated approaches to stereo and motion. Some of this work has concentrated on static scenes with variable lighting [4]. Others have focused on defining appropriate tempo-

ral integration windows, e.g., as part of correspondence [28] or simply reinforce disparity estimates from the previous frame using optical flow [9]. Further, combined stereo and motion estimation has been formulated in terms of PDEs [25, 11] as well as MRFs [27, 15]. Still other work has used direct methods for integrated recovery of structure and egomotion [24, 17]. The proposed research shares with these efforts an emphasis on tight integration of binocular imagery with time. It is novel in basing its matching on the representation of image spacetime in terms of local spatiotemporal orientation, which provides richer image descriptions than raw image intensities.

A major tool that is employed in the proposed approach is the representation of spacetime imagery in terms of oriented spatiotemporal structure. Various research documents optical flow recovery [2], tracking [3] and grouping [7] using spatiotemporal orientation tuned filters. More specifically, previous research has used the spatiotemporal quadric to capture orientation in image spacetime, with application to motion estimation, restoration and enhancement [10]. However, it appears none has exploited spatiotemporal orientation, in general, or the spatiotemporal quadric, specifically, for stereo disparity estimation. Previous stereo work has defined binocular correspondence based on a bank of spatial filters [13]. The proposed approach also extracts its measures of orientation via application of a filter bank; however, it is significantly different in employing filters that span both the space and time domains, thereby basing matching on a fundamentally richer representation.

In the light of previous research, the main contributions of this work are as follows. (i) The spatiotemporal quadric is proposed as a matching primitive for spatiotemporal stereo. This primitive captures both local spatial and temporal structure and thereby enables matching to account for both sources of data without need to estimate optical flow or 3D motion. (ii) The geometric relationships between corresponding spatiotemporal quadrics across binocular views are derived and used to motivate a match cost. The spatiotemporal match primitives and cost are incorporated in local and global matchers. (iii) Extensive empirical evaluation of these matchers is presented. Testing encompasses quantitative evaluation on laboratory acquired binocular video with ground truth and qualitative evaluation on more naturalistic imagery. The datasets and associated ground truth are available for download [23].

2. Technical Approach

2.1. Spatiotemporal matching primitive

In dealing with temporal sequences of binocular images, it is possible to conceptualize of stereo correspondence in terms of image spacetime, which naturally encompasses both spatial and temporal characteristics of local pattern structure, see Fig. 1a. While image spacetime can be op-

erated on directly, using intensities, consideration of local spatiotemporal orientation provides access to a richer representation. Local orientation has visual significance: orientations parallel to the image plane capture the spatial pattern of observed surfaces (e.g. spatial texture); whereas, orientations that extend into the temporal dimension capture dynamic aspects (e.g. motion). By integrating the temporal dimension into the primitive, matching will be inherently constrained to observe temporal coherence. Further, via combination of temporal and spatial structure in the descriptor, match ambiguities that might exist through consideration of only one data source have potential to be resolved.

To extract a representation of orientation from imagery, one can filter the data with oriented filters. In the current work, 3D Gaussian, second-derivative filters, G_2 , and their Hilbert transforms, H_2 [8], are applied to the data with responses pointwise rectified (squared) and summed. Filtering is executed across a set of 3D orientations given by unit column vectors, $\hat{\mathbf{w}}_i$. Hence, a measure of local energy, E , is computed according to

$$E(\mathbf{x}; \hat{\mathbf{w}}_i) = [G_2(\hat{\mathbf{w}}_i) * I(\mathbf{x})]^2 + [H_2(\hat{\mathbf{w}}_i) * I(\mathbf{x})]^2, \quad (1)$$

where $\mathbf{x} = (x, y, t)$ are spatiotemporal image coordinates, I is the image sequence and $*$ denotes convolution [8]. Filtering is applied separately to left and right image sequences. Here, filters are oriented along normals to icosahedron faces with antipodal directions identified, as this uniformly sample the sphere and spans 3D orientation for the employed filters. After filtering, every point in spacetime has an associated set of values that indicate how strongly oriented the local structure is along each spacetime direction.

To proceed, the individual energy measures are recast in terms of the spatiotemporal quadric. This particular representation captures local orientation as well as the variance of spacetime about that orientation. This construct captures the local shape of spacetime (e.g. point- vs. line- vs. plane-like) in addition to direction for a local descriptor that is richer than if (dominant) orientation alone is considered [10]. Furthermore, the quadric casts structure in terms of spacetime coordinates, $\mathbf{x} = (x, y, t)$, where it is convenient to formulate binocular match constraints. In the context of binocular matching, this quadric will be referred to as the **stequel**, spatio-temporal *quadric element*, Q . In particular,

$$Q = \sum_i \hat{E}_i \hat{\mathbf{w}}_i \hat{\mathbf{w}}_i^T, \quad (2)$$

where summation is across the set of filter orientations, $\hat{\mathbf{w}}_i$, and \hat{E}_i is the corresponding local energy response (1) normalized such that $\sum_i \hat{E}_i(\mathbf{x}) = 1$. In constructing Q , the dyadic product, $\hat{\mathbf{w}}_i \hat{\mathbf{w}}_i^T$, sets the local frame implied by orientation $\hat{\mathbf{w}}_i$ weighted by its response, \hat{E}_i , [10].

For a binocular sequence, the stequel, Q , is computed pointwise in spacetime and separately for the left and right image sequences to provide matching primitives; thus, it is

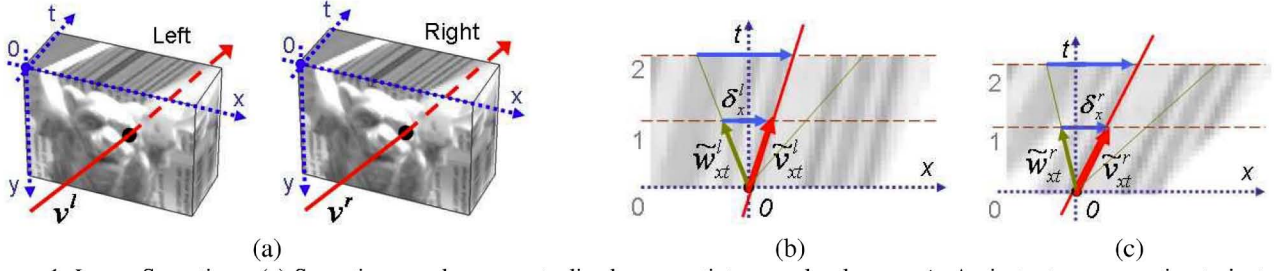


Figure 1. Image Spacetime. (a) Spacetime can be conceptualized as a spatiotemporal volume xyt . An instantaneous motion trajectory, v (shown in red), traces an orientation in this volume. (b) An exemplar xt slice of the spatiotemporal volume for the left view (c) The corresponding xt slice in the right view. \tilde{v}_{xt}^l and \tilde{v}_{xt}^r are the projections of the v^l and v^r onto the xt slice; w^l and w^r are arbitrary vectors (shown in green) in correspondence in xyt space and $\delta^r = \tilde{w}^r - \tilde{v}^r$, $\delta^l = \tilde{w}^l - \tilde{v}^l$ (shown in blue); $\delta^r = A\delta^l$ as explained in text.

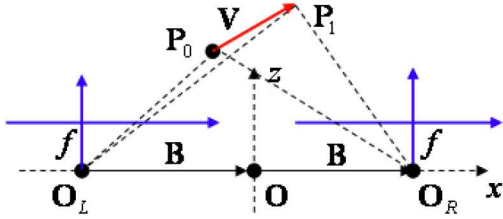


Figure 2. Stereo Geometry. A Euclidean coordinate system is set at the midpoint of the stereo baseline, O . Cameras are rectified with a half-baseline $\mathbf{B} = [B, 0, 0]^T$ and focal lengths f . Left and right optical centres are at $O^l = -\mathbf{B}$ and $O^r = \mathbf{B}$, resp. Point \mathbf{P} undergoes an arbitrary displacement \mathbf{V} from time 0 to 1.

parametrized as $Q^l(\mathbf{x})$ and $Q^r(\mathbf{x})$, in reference to the left and right views, resp. Significantly, the implied calculations are modest. The calculation of local energy is realized through steerable filters requiring only 3D separable convolution and pointwise nonlinearities and is thereby amenable to compact, efficient implementation [6]. Construction of Q from the filter responses requires only matrix summation, as specified in (2).

2.2. Spatiotemporal epipolar constraint

In establishing correspondence between binocular sequences, it is incorrect simply to seek the most similar stequels, as local spatiotemporal orientation is expected to change between views due to the geometry of the situation. In this section, constraint is derived between corresponding stequels subject to rectified and otherwise calibrated binocular viewing. This constraint is derived in two steps. First, the relationship between local spatiotemporal orientations in left and right image spacetime is derived as a 3D scene point \mathbf{P} suffers an arbitrary (infinitesimal) 3D displacement, \mathbf{V} , relative to the imaging system. While the relationship between left- and right-based flow has been investigated previously (e.g., [26]), the present derivation sets it in the light of left/right spatiotemporal orientation differences with application to disparity estimation. Second, the left/right flow relationships are generalized to capture the relationship between arbitrary orientations in left and right spacetimes. These results lead directly to the desired rela-

tionship between binocular stequels in correspondence.

In the following, bold and regular fonts denote vectors and scalars (resp.), uppercase denotes points relative to the world, lowercase denotes points relative to an image, superscripts l and r denote left and right cameras (resp.), subscripts x, y, z, t specify coordinate components, and vectors in image spacetime taken from time $t = 0$ to $t = 1$ will be distinguished further with tilde. As examples: $\mathbf{P}_t^l = [P_x^l \ P_y^l \ P_z^l]^T$ is the left camera representation of \mathbf{P} at time t ; $\mathbf{p}_t^l = [p_x^l \ p_y^l]^T$ is the left image coordinate of \mathbf{P}_t^l ; $\tilde{\mathbf{w}} = [w_x \ w_y \ 1]^T$ is a vector in image spacetime xyt from $t = 0$ to $t = 1$.

Left-Right Flow Relationship. Consider how a 3D point, \mathbf{P} , is observed by the cameras as a function of time, t , while it is displaced along 3D direction, \mathbf{V} . The geometry of the situation is shown in Fig. 2. Cameras share a common intrinsic matrix $\mathbf{K} = \begin{bmatrix} f & 0 & 0 \\ 0 & f & 0 \\ 0 & 0 & 1 \end{bmatrix}$, where other components of the matrix are accounted for by calibration and neglected. At time t , the projections of \mathbf{P} to the left and right views are given by

$$\mathbf{P}_t^l = \mathbf{K}((\mathbf{P}_{t=0} - \mathbf{B}) + t\mathbf{V}) = \mathbf{P}_{t=0}^l + t\mathbf{K}\mathbf{V} \quad (3)$$

$$\mathbf{P}_t^r = \mathbf{K}((\mathbf{P}_{t=0} + \mathbf{B}) + t\mathbf{V}) = \mathbf{P}_{t=0}^r + t\mathbf{K}\mathbf{V}.$$

Note that both moving and stationary points are encompassed in this formulation, as \mathbf{V} is arbitrary. The corresponding image coordinates are found in the usual way, e.g. for the left view

$$\mathbf{p}^l = \begin{bmatrix} p_x^l \\ p_y^l \end{bmatrix} = \begin{bmatrix} P_x^l/P_z^l \\ P_y^l/P_z^l \end{bmatrix} = \frac{1}{P_z^l} \begin{bmatrix} P_x^l \\ P_y^l \end{bmatrix} = Z^{-1}\mathbf{P}_{2 \times 1}^l, \quad (4)$$

where $P_z^l = Z$ is the distance along the Z -axis to the point of regard, \mathbf{P} , and $\mathbf{P}_{2 \times 1}$ is the upper 2×1 component of \mathbf{P} . Analogously for right view, $\mathbf{p}^r = Z^{-1}\mathbf{P}_{2 \times 1}^r$.

In the image spacetime coordinate system, xyt , without loss of generality, consider flows \tilde{v}^l and \tilde{v}^r in the left and right views from temporal instance 0 to 1:

$$\tilde{\mathbf{v}}^l = \begin{bmatrix} \mathbf{P}_{t=1}^l - \mathbf{P}_{t=0}^l \\ v_t^l \end{bmatrix} = \begin{bmatrix} \mathbf{P}_{t=1}^l - \mathbf{P}_{t=0}^l \\ 1 \end{bmatrix}, \quad (5)$$

where $v_t^l = 1$ by definition, as time has been taken from

$l = 0$ to $l = 1$. Analogously for the right view

$$\tilde{\mathbf{v}}^r = \begin{bmatrix} \mathbf{P}_{t-1}^r - \mathbf{P}_{t-0}^r \\ 1 \end{bmatrix}. \quad (6)$$

To relate the left and right spatiotemporal orientations, it is useful to cast the left-camera flow vectors (5) and their right camera counterparts in terms of temporally varying position (3) and (4). Left camera-based flow is given by (5) and substitution from (4) yields

$$\tilde{\mathbf{v}}_{2 \times 1}^l = Z_{t=1}^{-1} \mathbf{P}_{2 \times 1, t=1}^l - Z_{t=0}^{-1} \mathbf{P}_{2 \times 1, t=0}^l.$$

Further substitution for \mathbf{P}^l according to (3) and letting all subscripts pertain to time (i.e. 0 and 1 denote $t = 0$ and $t = 1$, resp.) yields

$$\tilde{\mathbf{v}}_{2 \times 1}^l = \frac{Z_0 - Z_1}{Z_0 Z_1} \bar{\mathbf{K}} \mathbf{P}_0 + \frac{1}{Z_1} \bar{\mathbf{K}} \mathbf{V} - \frac{Z_0 - Z_1}{Z_0 Z_1} \bar{\mathbf{K}} \mathbf{B}, \quad (7)$$

where $\bar{\mathbf{K}} = \mathbf{K}_{2 \times 3}$ is the top two rows of \mathbf{K} . Similarly, for the right camera-based flow

$$\tilde{\mathbf{v}}_{2 \times 1}^r = \frac{Z_0 - Z_1}{Z_0 Z_1} \bar{\mathbf{K}} \mathbf{P}_0 + \frac{1}{Z_1} \bar{\mathbf{K}} \mathbf{V} + \frac{Z_0 - Z_1}{Z_0 Z_1} \bar{\mathbf{K}} \mathbf{B}. \quad (8)$$

Finally, the relationship between the left (7) and right (8) flows is revealed by taking their difference

$$\tilde{\mathbf{v}}^r - \tilde{\mathbf{v}}^l = \begin{bmatrix} 2(Z_0 - Z_1) \bar{\mathbf{K}} \mathbf{B} / (Z_0 Z_1) \\ 0 \end{bmatrix} = \begin{bmatrix} \Delta \\ 0 \\ 0 \end{bmatrix}, \quad (9)$$

where $\Delta = 2Bf(Z_0 - Z_1) / (Z_0 Z_1)$ captures the instantaneous change in disparity.

General Left/Right Orientation Relationship. The relationship (9) was derived only for dominant motion orientation; whereas, stequels capture information from *all* directions $\tilde{\mathbf{w}}$ in (x, y, t) , which now are considered.

Consider directions $\tilde{\mathbf{w}}^r$ and $\tilde{\mathbf{w}}^l$ in the left and right views, resp., that are in binocular correspondence, but otherwise arbitrary in (x, y, t) . Discounting the effects of right and left flows, $\tilde{\mathbf{v}}^r$ and $\tilde{\mathbf{v}}^l$, yields vectors

$$\delta^r = \tilde{\mathbf{w}}^r - \tilde{\mathbf{v}}^r = \begin{bmatrix} \delta_x^r & \delta_y^r & 0 \end{bmatrix}^\top, \quad (10)$$

$$\delta^l = \tilde{\mathbf{w}}^l - \tilde{\mathbf{v}}^l = \begin{bmatrix} \delta_x^l & \delta_y^l & 0 \end{bmatrix}^\top, \quad (11)$$

that capture the purely spatial orientation of corresponding elements (see Fig. 1b,c). For the special case of fronto-parallel surfaces $\delta^r = \delta^l$, i.e. disregarding motion, oriented texture appears the same across binocular views. For the more general case where surfaces are slanted with respect to the imaging system, the imaged orientation of corresponding elements changes across views, even in the absence of motion. For present matters, this change can be modeled by a linear transformation $\delta^r = \mathbf{A} \delta^l$. Considering that the third element of the δ vectors is always zero by construction, and $\delta_y^r = \delta_y^l$ due to conventional stereo epipolar constraints for rectified setups, this relationship takes the form

$$\delta^r = \mathbf{A} \delta^l, \text{ where } \mathbf{A} = \begin{bmatrix} a_1 & a_2 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}. \quad (12)$$

Substituting (10), (11) into (12) and rearranging yields,

$$\tilde{\mathbf{w}}^r = \mathbf{A} \tilde{\mathbf{w}}^l - \mathbf{A} \tilde{\mathbf{v}}^l + \tilde{\mathbf{v}}^r. \quad (13)$$

Further substitution of (9) results in

$$\begin{aligned} \tilde{\mathbf{w}}^r &= \mathbf{A} \tilde{\mathbf{w}}^l + \left(-\mathbf{A} \tilde{\mathbf{v}}^l + \tilde{\mathbf{v}}^l + \begin{bmatrix} \Delta & 0 & 0 \end{bmatrix}^\top \right) \\ &= \begin{bmatrix} a_1 & a_2 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \tilde{\mathbf{w}}^l + \begin{bmatrix} 1 - a_1 & -a_2 & \Delta \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \tilde{\mathbf{v}}^l \\ &= \begin{bmatrix} a_1 & a_2 & ((1 - a_1)\tilde{v}_x^l - a_2\tilde{v}_y^l + \Delta) \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \tilde{\mathbf{w}}^l \end{aligned} \quad (14)$$

Finally, letting $h_1 = a_1 - 1$, $h_2 = a_2$ and $h_3 = ((1 - a_1)\tilde{v}_x^l - a_2\tilde{v}_y^l + \Delta)$ yields the desired transformation between arbitrary corresponding vectors $\tilde{\mathbf{w}}^l$ and $\tilde{\mathbf{w}}^r$

$$\tilde{\mathbf{w}}^r = \mathbf{H} \tilde{\mathbf{w}}^l, \text{ where } \mathbf{H} = \begin{bmatrix} 1 + h_1 & h_2 & h_3 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad (15)$$

With (15) in place, it is possible to relate corresponding stequels. By design, (2), stequel \mathbf{Q} reveals the amount of intensity variation along all directions in spacetime, and the response ϕ to unit direction $\hat{\mathbf{w}} = \mathbf{w} / \sqrt{\mathbf{w}^\top \mathbf{w}}$ is

$$\phi = \hat{\mathbf{w}}^\top \mathbf{Q} \hat{\mathbf{w}}, \quad (16)$$

see, e.g., [10]. Assuming that spatiotemporal correspondences vary in orientation pattern, but not in the intensity per se¹, the responses, ϕ^l, ϕ^r , of corresponding stequels, $\mathbf{Q}^l, \mathbf{Q}^r$, must be the same for related directions, $\hat{\mathbf{w}}^l, \hat{\mathbf{w}}^r$, i.e.

$$\hat{\mathbf{w}}^{l\top} \mathbf{Q}^l \hat{\mathbf{w}}^l = \hat{\mathbf{w}}^{r\top} \mathbf{Q}^r \hat{\mathbf{w}}^r.$$

Expanding the normalizations of $\hat{\mathbf{w}}^l$ and $\hat{\mathbf{w}}^r$ and substituting from (15) produces

$$\frac{\tilde{\mathbf{w}}^{l\top} \mathbf{Q}^l \tilde{\mathbf{w}}^l}{\tilde{\mathbf{w}}^{l\top} \tilde{\mathbf{w}}^l} = \frac{\tilde{\mathbf{w}}^{l\top} \mathbf{H}^\top \mathbf{Q}^r \mathbf{H} \tilde{\mathbf{w}}^l}{\tilde{\mathbf{w}}^{l\top} \mathbf{H}^\top \mathbf{H} \tilde{\mathbf{w}}^l},$$

while noticing that $\tilde{\mathbf{w}}^l = \|\tilde{\mathbf{w}}^l\| \hat{\mathbf{w}}^l$ yields

$$\frac{\hat{\mathbf{w}}^{l\top} \mathbf{Q}^l \hat{\mathbf{w}}^l}{\hat{\mathbf{w}}^l \hat{\mathbf{w}}^l} = \frac{\hat{\mathbf{w}}^{l\top} \mathbf{H}^\top \mathbf{Q}^r \mathbf{H} \hat{\mathbf{w}}^l}{\hat{\mathbf{w}}^{l\top} \mathbf{H}^\top \mathbf{H} \hat{\mathbf{w}}^l}. \quad (17)$$

Since (17) holds for arbitrary orientations $\hat{\mathbf{w}}^l$ when \mathbf{Q}^l and \mathbf{Q}^r are stequels in correspondence, it provides the sought for general constraint on binocular stequels. It will be referred to as the *stequel correspondence constraint* and used to derive an approach to stereo matching.

2.3. Stequel match cost

To determine whether two stequels $\mathbf{Q}^l(x, y, t)$ and $\mathbf{Q}^r(x + d, y, t)$ are in correspondence with disparity d , a match cost must be defined. In this section, this cost is derived based on the stequel correspondence constraint, (17), and is taken as the error residual that results from solving for $\mathbf{h} = [h_1 \ h_2 \ h_3]^\top$ given two candidate stequels.

For a given direction vector $\hat{\mathbf{w}}_m^l$ at some particular orientation m and matching stequels, \mathbf{Q}^l and \mathbf{Q}^r , the stequel correspondence constraint, (17), yields a quadratic equation

¹This is a weak form of brightness constancy as any additive and multiplicative intensity offsets between correspondences are compensated for by the bandpass and normalized filters used in stequel construction (2).

in the unknowns of \mathbf{h} of the form

$$f_m(\mathbf{h}) = (\hat{\mathbf{w}}_m^{l\top} \mathbf{Q}^l \hat{\mathbf{w}}_m^l) (\hat{\mathbf{w}}_m^{l\top} \mathbf{H}^\top \mathbf{H} \hat{\mathbf{w}}_m^l) - (\hat{\mathbf{w}}_m^{l\top} \hat{\mathbf{w}}_m^l) (\hat{\mathbf{w}}_m^{l\top} \mathbf{H}^\top \mathbf{Q}^r \mathbf{H} \hat{\mathbf{w}}_m^l) = 0. \quad (18)$$

Taking a set of M directions, reasonably selected along the same spanning set of directions used to construct \mathbf{Q}^l , yields a set of M equations in the three unknowns of \mathbf{h} . Thus, \mathbf{h} can be estimated by minimizing a sum of squared errors

$$E_A = \sum_{m=1}^M f_m(\mathbf{h})^2, \quad (19)$$

which is quartic in the entries of \mathbf{h} . While such a solution could be sought through analytic or numerical means, it has potential to be expensive to compute and noise sensitive owing to its order. Therefore, it is useful to linearize each error Eqn. (18) through expansion as a Taylor series in \mathbf{h} and retention of terms only through first-order to get

$$g_m(\mathbf{h}) = f_m(\mathbf{0}) + \nabla f_m^l(\mathbf{0})\mathbf{h}, \quad (20)$$

with $\mathbf{0}$ being the $M \times 1$ zero vector. Using (20), the final function to be minimized with respect to \mathbf{h} becomes

$$E_2 = \sum_{m=1}^M (f_m(\mathbf{0}) + \nabla f_m^l(\mathbf{0})\mathbf{h})^2, \quad (21)$$

which is simply quadratic in the elements of \mathbf{h} , and thereby can be solved for via standard linear least-squares. More specifically, letting $\mathbf{G} = [\nabla f_1^l(\mathbf{0}), \nabla f_2^l(\mathbf{0}), \dots, \nabla f_M^l(\mathbf{0})]^\top$ and $\mathbf{c} = -[f_1(\mathbf{0}), f_2(\mathbf{0}), \dots, f_M(\mathbf{0})]^\top$ yields

$$\mathbf{h} = (\mathbf{G}^\top \mathbf{G})^{-1} \mathbf{G}^\top \mathbf{c}; \quad (22)$$

$$E_2 = \|\mathbf{G}\mathbf{h} - \mathbf{c}\|_2^2 = \mathbf{c}^\top \mathbf{c} - (\mathbf{G}^\top \mathbf{c})^\top (\mathbf{G}^\top \mathbf{G})^{-1} \mathbf{G}^\top \mathbf{c}.$$

For two stequels under consideration for stereo correspondence this residual, E_2 , will serve as the local match cost².

3. Empirical Evaluation

A software implementation has been developed that inputs a binocular video, computes stequels $\mathbf{Q}^l(x, y, t)$ and $\mathbf{Q}^r(x, y, t)$ for both sequences according to formula (2)³ and calculates the local match cost, (22), for any given disparity d , i.e., for stequels related as $\mathbf{Q}^l(x, y, t)$ and $\mathbf{Q}^r(x + d, y, t)$. To show the applicability of this approach to disparity estimation, the local match cost, (22), has been embedded in a coarse-to-fine local block-matching algorithm with shiftable windows [22] working over a Gaussian pyramid and also in a global graph-cuts with occlusions matcher [14] operating at the finest scale only; these matchers will be

²Significantly, preliminary experiments showed that match cost based on the linearized error, (21), yielded slightly superior results to considering the original nonlinear error, (19). This can be explained by noting that interest is in a discriminative error measure that reliably penalizes bad matches, and not in the precise error value per se (see [23] for details).

³Preliminary experiments considered alternative stequel definition $\mathbf{Q} = \sum \nabla I (\nabla I)^\top$, with $\nabla = [I_x I_y I_t]^\top$ is the spatiotemporal gradient and summation is over local spacetime regions (see [23] for details), but found generally inferior results; so, this approach is considered no further.

denoted **ST-local** and **ST-global**. Pixel-based disparity estimates are brought to subpixel precision via a Lucas-Kanade type refinement for stequels [1, 23].

To compare with non-stequel matching, versions of the local and global matchers that work simply on single left/right frame pixel comparisons are considered; these matchers will be denoted **noST-local** and **noST-global**, resp. Here, the normalized cross-correlation was used as the data cost term for local and global matching. Finally, to compare to an alternative method for enforcing temporal coherence, optical flow is estimated and used to define a spatiotemporal direction for match cost aggregation that operates over an equivalent number of frames as does the oriented filtering used in stequel construction (1). Here, optical flow is recovered from the stequel representation itself (see [23] and [10] for discussion) to make the comparison fair. The optical flow-based temporal aggregation is used only in conjunction with the local matcher, as incorporation into the global matcher by constructing a spatiotemporal MRF graph [15] is beyond the scope of this paper. The local flow-based aggregation matcher will be denoted **flowAg-local**.

Three data sets are considered. The first is a laboratory sequence (*Lab1*) captured with BumbleBee stereo camera [19] with (framewise) ground truth disparity and discontinuity maps recovered according to a well-known structured light approach [21], see Fig. 3. This scene includes planes slanted in depth with texture oriented along epipolar lines (upper-central part of the scene), various bar-plane arrangement with identical repetitive textures (lower-central part of the scene) and complicated objects with non-trivial 3D boundaries and non-Lambertian materials (e.g., the teddy bear and gargoyle). For this sequence the stereo camera makes a complicated motion that translates along horizontal and depth axes, while part of the scene moves up and down; both camera and scene are on motorized stages.

Visual inspection of the image results (Fig. 3) shows that **noST-local** performs relatively poorly. Planar regions with epipolar aligned texture are generally difficult. Simple temporal aggregation provided by **flowAg-local** is seen to improve on these difficulties; however, performance degrades near 3D boundaries due to unreliable recovery of flow estimates in such areas. **ST-local** does the best of the three local matchers as its *ability to include temporal information allows it to resolve match ambiguities without explicit flow recovery*. As particular improvements of **ST-local** over **noST-local** and **flowAg-local**, consider the lower right and left regions marked with red rectangles in Fig. 3, which highlight the complex outline of the gargoyle wings and the vertical bar in front of plane both having identical textures (camouflage). **ST-local** is quite accurate in these challenging regions, while the other local methods perform relatively poorly. Objects located at different depths in space give rise

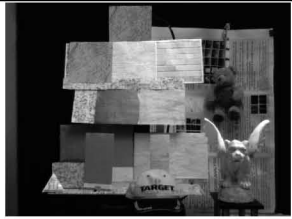


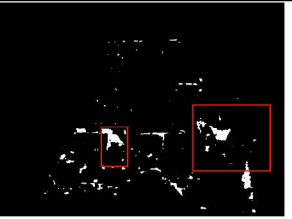



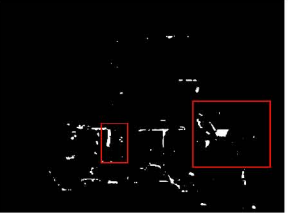
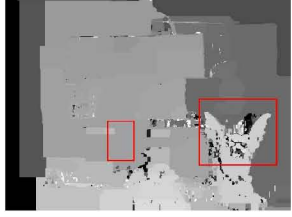



Lab 1 Left frame 12	GT disparity	flowAg-local disparity	flowAg-local error
			
noST-local disparity	noST-local error	ST-local disparity	ST-local error
			
noST-global disparity	noST-global error	ST-global disparity	ST-global error
			

Figure 3. *Lab1* Tests. Example left frame 12 (out of 28 frames) with ground truth disparity at a single time instance. Labeled boxes show recovered disparity maps for compared algorithms and disparity-ground truth absolute differences. A few regions of particular interest in comparing results are highlighted with red rectangles, best seen in color.





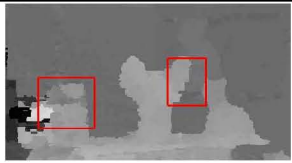




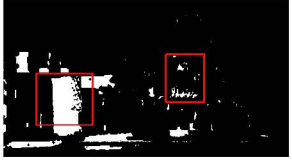


Lab 2 Left frame 10	GT disparity	flowAg-local disparity	flowAg-local error
			
noST-local disparity	noST-local error	ST-local disparity	ST-local error
			
noST-global disparity	noST-global error	ST-global disparity	ST-global error
			

Figure 4. *Lab2* Tests. Example left frame 10 (out of 40 frames) with ground truth disparity at a single time instance. Labeled boxes show recovered disparity maps for compared algorithms and disparity-ground truth absolute differences. A few regions of particular interest in comparing results are highlighted with red rectangles, best seen in color.

to different image motions, even if they undergo the same world motion – and this difference is captured with stequels not allowing for improper matches.

For the global matchers, it is seen even with **noST-global** that it is possible to recover more precisely the complicated 3D boundaries and to achieve good disparity estimates

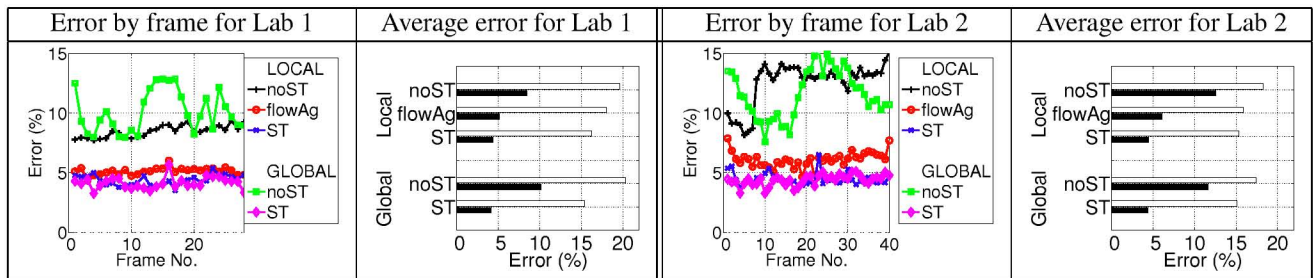


Figure 5. Error statistics for the *Lab1* and *Lab2* Tests. An error is taken as greater than 1 pixel discrepancy between estimated and groundtruth disparity. Bar plots show average error across entire sequences: White bars are for points within 5 pixels of a surface discontinuity; black bars show overall error. Error by frame plots show percentage of points in error overall for each frame separately.

in low texture regions via propagation from better defined boundary matches. However, **noST-global** performs poorly in the regions with epipolar aligned texture and camouflage, as initially incorrect estimates are not subsequently corrected. While increasing the smoothness improves on epipolar-aligned textures, it comes at the expense of camouflage resolution and vice versa. In comparison, **ST-global** is able to recover disparity reliably in these regions, as *once again the stequel representation supports proper resolution of situations that are ambiguous from the purely spatial information*. Another apparent advantage of the **ST-global** is more temporally consistent results – occasional mismatches in **noST-global** can be significantly amplified by propagating into nearby regions.

A second lab sequence, *Lab2*, is constructed in the same controlled environment as *Lab1*, but acquired with significant depth motion and out-of-plane rotation. This particular motion configuration is the most difficult for spatiotemporal stereo, as it results in significantly different left and right spatiotemporal volumes due to slanted surfaces and depth motion. Furthermore, large image motions are present in the individual left and right sequences. Figure 4 presents sample frame results for all five algorithmic instantiations considered above. Here, the conclusions reached from the analysis of *Lab1* are reinforced. With respect to the local methods, **ST-local** provides the most benefit both in weakly textured regions and near 3D boundaries. The performance of **flowAg-local** is hampered by large image motions, which are problematic to recover explicitly in this case; whereas, *direct stequel-based matching is still able to capitalize on temporal information without resolving flow and thereby operates well in the presence of nontrivial motions*. With respect to the global methods, the stequel-based matching **ST-global** significantly outperforms its pixel-based counterpart **noST-global**, especially for weakly-textured highly slanted foreground surfaces.

Error plots for both *Lab1* and *Lab2* quantify the improvements of stequel-based matching in comparison to rivals **noST** and **flowAg** (Fig. 5). Average errors across the sequences show the benefit of stequels near discontinuities and overall for both local and global matchers. Plots of er-

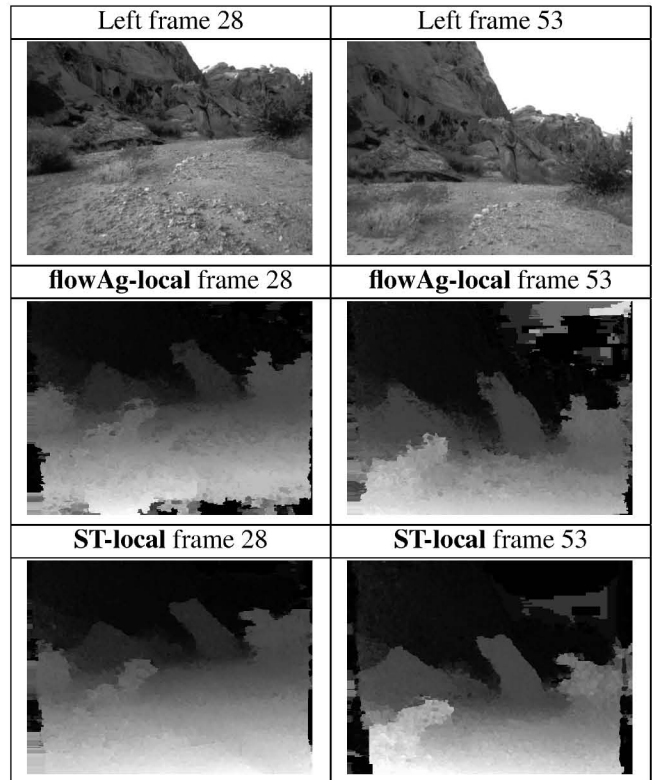


Figure 6. *Rover Tests*. Top row shows left view at frames 28 and 53. Recovered disparity maps at corresponding times are shown below for two algorithms under consideration.

ror/frame reinforce the average improvements, but also document improved temporal coherence, as the stequel-based plots vary relatively little across frames, especially in comparison to purely spatial matching provided by **noST**. Incorporation of the temporal dimension also benefits **flowAg**, as its frame-by-frame statistics are relatively stable (albeit overall inferior to stequels); however, the more naturalistic imagery of the next example further emphasizes the superior temporal coherence offered by stequels, even in comparison to **flowAg**.

The third data set, *Rover*, is an outdoor sequence acquired from a robot rover traversing rugged terrain, including a receding foreground plane, a central diagonal rock

outcropping, left side cliff, various boulders and bushes. Here the comparison focuses on the improvements to temporal coherence offered by **ST-local** over the rival method for consideration of temporal information, **flowAg-local**. As results of depicted frames show, flow-based aggregation, while providing mostly temporally coherent estimates is inferior at recovery of 3D boundaries (boulders' outlines) and still susceptible to occasional gross errors (*e.g.* on the ground plane) due to errors in the recovered flow. In comparison, stequel-based matching, **ST-local**, does not exhibit such problems, as it uses spatiotemporal information in a more direct and complete way.

4. Discussion

This paper described a novel approach to recovering temporally coherent disparity estimates using stequels as a spatiotemporal matching primitive. Temporal coherence arises naturally, as the primitives and the match cost inherently involve the temporal dimension. Further, matches that are ambiguous when considering only spatial pattern are resolved through the inclusion of temporal information. The stequel matching machinery is simple and involves linear computations only, (22). Thorough experimental evaluation on various datasets shows the benefit of stequel matching as incorporated both in local and global algorithms. Stereo sequences with ground truth have been introduced and are available online for comparison with other algorithms, [23].

A particularly notable benefit of stequel matching is the ability to incorporate temporal information *without* image motion recovery. Optical flow estimation is challenging near 3D boundaries, weakly-textured regions and susceptible to an aperture problem – importantly, this paper demonstrated that stequels are powerful in exactly these situations and provide truly temporally coherent estimates with fewer isolated gross errors. Apparently, stequels allow stereo matching to capitalize on available spatiotemporal structure, even when optical flow recovery is difficult. Further, note that it is non-trivial to model continuity in time with, *e.g.* an MRF prior model as, strictly speaking, temporal graph links have to be defined by flow (as in [15]). Stequels, on the other hand, are directly applicable to standard 2D MRF graphs and their successful performance has been documented in this paper.

In conclusion, a computationally tractable and simple solution to spatiotemporal stereo has been presented, which proved to be very reliable, versatile and robust in practice. Future work will concentrate on exploiting the spatiotemporal profile for explicit non-Lambertian and multi-layer matching, as well as extensions to 3D motion recovery.

Acknowledgements

This work was supported by NSERC and MDA Space Missions. P. Jasiobedski and S. Se provided the *Rover* sequence.

References

- [1] S. Baker and I. Matthews. Lucas-Kanade 20 years on: A unifying framework. *IJCV*, 56(3):221–255, 2004.
- [2] S. Beauchemin and J. Barron. The computation of optical flow. *ACM Comp. Surv.*, 27:433–467, 1995.
- [3] K. Cannons and R. P. Wildes. Spatiotemporal oriented energy features for visual tracking. In *ACCV*, pages 532–543, 2007.
- [4] J. Davis, R. Ramamoorthi, and S. Rusinkiewicz. Spacetime stereo: A unifying framework for depth from triangulation. *TPAMI*, 27(2):296–302, 2005.
- [5] D. Demirdjian and T. Darrell. Using multiple-hypothesis disparity maps and image velocity for 3D motion estimation. *IJCV*, 47:219–228, 2002.
- [6] K. G. Derpanis and J. Gryn. 3-D n-th derivative of Gaussian separable steerable filters. *ICIP*, 3:553–556, 2005.
- [7] K. G. Derpanis and R. P. Wildes. Early spatiotemporal grouping with a distributed oriented energy representation. In *CVPR*, 2009.
- [8] W. T. Freeman and E. H. Adelson. The design and use of steerable filters. *TPAMI*, 13(9):891–906, 1991.
- [9] M. Gong. Enforcing temporal consistency in real-time stereo estimation. In *ECCV*, pages 564–577, 2006.
- [10] G. Granlund and H. Knutsson. *Signal Processing for Computer Vision*. Kluwer, 1995.
- [11] F. Huguet and F. Devernay. A variational method for scene flow estimation from stereo sequences. In *ICCV*, 2007.
- [12] M. Jenkin and J. K. Tsotsos. Applying temporal constraints to the dynamic stereo problem. *CVGIP*, 33:16–32, 1986.
- [13] D. G. Jones and J. Malik. A computational framework for determining stereo correspondence from a set of linear spatial filters. In *ECCV*, pages 395–410, 1992.
- [14] V. Kolmogorov and R. Zabih. Computing visual correspondence with occlusions using graph cuts. In *ICCV*, pages 508–515, 2001.
- [15] E. S. Larsen, P. Mordohai, M. Pollefeys, and H. Fuchs. Temporally consistent reconstruction from multiple video streams. In *ICCV*, pages 1–8, 2007.
- [16] S. Malassiotis and M. G. Strintzis. Model-based joint motion and structure estimation from stereo images. *CVIU*, 65(1):79–94, 1997.
- [17] R. Mandelbaum, G. Salgian, and H. Sawhney. Correlation-based estimation of ego-motion and structure from motion and stereo. In *ICCV*, pages 544–550, 1999.
- [18] J. Neumann and Y. Aloimonos. Spatio-temporal stereo using multi-resolution subdivision surfaces. *IJCV*, 47(1-3):181–193, 2002.
- [19] Point Grey Research. <http://www.ptgrey.com>, 2008.
- [20] J.-P. Pons, R. Keriven, O. D. Faugeras, and G. Hermosillo. Variational stereovision and 3D scene flow estimation with statistical similarity measures. In *ICCV*, pages 597–602, 2003.
- [21] D. Scharstein and R. Szeliski. High-accuracy stereo depth maps using structured light. *CVPR*, 1:195–202, 2003.
- [22] M. Sizintsev and R. P. Wildes. Efficient stereo with accurate 3-D boundaries. In *BMVC*, pages 237–246, 2006.
- [23] M. Sizintsev and R. P. Wildes. Spatiotemporal stereo via spatiotemporal quadric element (stequel) matching. Technical Report CS-2008-04, York University, 2008.
- [24] G. P. Stein and A. Shashua. Direct estimation of motion and extended scene structure from a moving stereo rig. *CVPR*, 1:211–218, 1998.
- [25] C. Strecha and L. van Gool. Motion-stereo integration for depth estimation. In *ECCV*, pages 170–185, 2002.
- [26] A. M. Waxman and J. H. Duncan. Binocular image flows: Steps toward stereo-motion fusion. *TPAMI*, 8(6):715–729, 1986.
- [27] O. Williams, M. Isard, and J. MacCormick. Estimating disparity and occlusions in stereo video sequences. *CVPR*, pages 250–257, 2005.
- [28] L. Zhang, B. Curless, and S. M. Seitz. Spacetime stereo: Shape recovery for dynamic scenes. In *CVPR*, pages 367–374, 2003.
- [29] Z. Zhang and O. D. Faugeras. Three-dimensional motion computation and object segmentation in a long sequence of stereo frames. *IJCV*, 7(3):211–241, 1992.