

Feature Correspondence Via Graph Matching: Models and Global Optimization

Lorenzo Torresani¹, Vladimir Kolmogorov², and Carsten Rother¹

¹ Microsoft Research Ltd., Cambridge, UK
{ltorre, carrot}@microsoft.com

² University College London, UK
vnk@adastral.ucl.ac.uk

Abstract. In this paper we present a new approach for establishing correspondences between sparse image features related by an unknown non-rigid mapping and corrupted by clutter and occlusion, such as points extracted from a pair of images containing a human figure in distinct poses. We formulate this matching task as an energy minimization problem by defining a complex objective function of the appearance and the spatial arrangement of the features. Optimization of this energy is an instance of graph matching, which is in general a NP-hard problem. We describe a novel graph matching optimization technique, which we refer to as dual decomposition (DD), and demonstrate on a variety of examples that this method outperforms existing graph matching algorithms. In the majority of our examples DD is able to find the global minimum within a minute. The ability to globally optimize the objective allows us to accurately learn the parameters of our matching model from training examples. We show on several matching tasks that our learned model yields results superior to those of state-of-the-art methods.

1 Introduction

Feature correspondence is one of the fundamental problems of computer vision and is a key ingredient in a wide range of applications including object recognition, 3D reconstruction, mosaicing, motion segmentation, and image morphing. Several robust algorithms (see e.g. [1,2]) exist for registration of images of static scenes and for visual correspondence under rigid motion. These methods typically exploit powerful constraints (e.g. epipolar constraints) to reduce the search space and disambiguate the correspondence problem. However, such constraints do not apply in the case of non-rigid motion or when matching different object instances. A popular approach in these cases is to discard the information about the spatial layout of features, and to find correspondences using appearance only. For example, many object recognition methods (see e.g. [3,4]) represent images as orderless sets of local appearance descriptors, known as bags of features. Recent work [5] has suggested that for many correspondence problems, learned appearance-based models perform similarly or better than state-of-the-art structural models exploiting information about spatial arrangement of features. This is primarily due to the challenges posed by the optimization and training of structural models, which often require approximate solution of NP-hard problems. In this paper we contrast this theory, and demonstrate that a complex structural model for image matching

can be learned and optimized successfully. We cast the visual correspondence problem as an energy minimization task by defining a complex image matching objective depending on (i) feature appearance, (ii) geometric compatibility of correspondences, and (iii) spatial coherence of matched features. Additionally, we impose a uniqueness constraint allowing at most one match per feature. We introduce a novel algorithm to minimize this function based on the dual decomposition approach (DD) from combinatorial optimization, see e.g. [6,7,8,9,10,11]. The DD method works by maximizing a lower bound on the energy function. The value of the lower bound can be used to gauge the distance from the global minimum and to decide when to stop the optimization, in the event the global minimum cannot be found. For the majority of our examples DD finds the global minimum in reasonable time, and otherwise provides a solution whose cost is very close to the optimum. In contrast, previously proposed optimization methods such as [12,13] often fail to compute good solutions for our energy function. Our experimental evaluation shows that the model and the algorithm presented in this paper can be applied to a wide range of image matching problems with results matching or exceeding those of existing algorithms [5,14].

1.1 Relation to Previous Work

Models for feature matching. Our technique is loosely related to algorithms that find correspondences by matching appearance descriptors under smooth spatial transformations (see e.g. [15,16]). However, unlike such approaches, our method does not make a parametric assumption about the transformation relating the input images, and thus can be used in a wider range of applications. Belongie et al. [14] inject spatial smoothness in the match by means of an iterative technique that alternates between finding correspondences using shape features, and computing a regularized transformation aligning the matching features. The shape descriptors are recomputed in each iteration after the warping. Since the objective is changed at each iteration, the convergence properties of this algorithm are not clear. Our approach is most closely related to the work of Berg et al. [17], and Leordeanu and Hebert [18], who formulate visual correspondence as a graph matching problem by defining an objective including terms based on appearance similarity as well as geometric compatibility between pairs of correspondences. Our model differs from those in [17,18] in several ways. The methods proposed in [17] and [18] handle outliers by removing low-confidence correspondences from the obtained solutions. Instead, we include in our energy an explicit occlusion cost, as for example previously done in [19]. Thus our algorithm solves for the outliers as part of the optimization. We add to the objective a spatial coherence term, favoring spatial aggregation of matched features, which reduces the correspondence error on our examples. We also show that geometric penalty functions defined in local neighborhoods provide more accurate correspondences than global geometric costs, such as those used in [17] and [18]. Finally, we use the method of Liu et al. [20] to learn the parameter values for the model from examples, thus avoiding the need of manual parameter tuning.

Graph matching optimization. Graph matching is a challenging optimization problem which received considerable attention in the literature (see [21] for a comprehensive survey of methods). Proposed techniques include the graduated assignment algorithm of

Gold and Rangarajan [22], spectral relaxation methods [18,12], COMPOSE method of Duchi et al. [13], Maciel and Costeira [23] reduce the problem to concave minimization and apply the exact method in [24]. Torr [19] and Schellewald and Schnörr [25] use semi-definite programming (SDP) relaxation for graph matching. Among these papers, only [23] and [25] report obtaining optimal (or near optimal) solutions. The method in [23] was tested only on a single example with quadratic costs. We conjecture that on practical challenging instances this method will suffer from an exponential explosion¹. As shown in [12], the SDP relaxation approach in [25] scales quite poorly and is too expensive for problems of reasonable size.

2 Energy Function

We now describe the energy function of our matching model. Let P' and P'' be the sets of features extracted from the two input images. We denote with $A \subseteq P' \times P''$ the set of potential assignments between features in the two sets. We will use the terms assignment and correspondence interchangeably to indicate elements of A . We represent a *matching configuration* between the two point sets as a binary valued vector $\mathbf{x} \in \{0, 1\}^A$. Each correspondence $a \in A$ indexes an entry x_a in the vector \mathbf{x} . A correspondence a is active if $x_a = 1$, and it is inactive otherwise. We define an energy function $E(\mathbf{x})$ modeling our matching problem assumptions. This will allow us to formulate the matching task as minimization of $E(\mathbf{x})$. In this paper we consider matching problems where at most one active correspondence per feature is allowed. This requirement is known as the uniqueness constraint and it is commonly used in correspondence problems. In order to enforce this condition we define the constraint set M :

$$M = \{\mathbf{x} \in \{0, 1\}^A \mid \sum_{a \in A(p)} x_a \leq 1 \quad \forall p \in P\} \quad (1)$$

where $P = P' \cup P''$ is the set of features from both images, and $A(p)$ is the set of correspondences involving feature p . The goal is to find the configuration $\mathbf{x} \in M$ minimizing $E(\mathbf{x})$. We define our energy as a weighted sum of four energy terms:

$$E(\mathbf{x}) = \lambda^{\text{app}} E^{\text{app}}(\mathbf{x}) + \lambda^{\text{occl}} E^{\text{occl}}(\mathbf{x}) + \lambda^{\text{geom}} E^{\text{geom}}(\mathbf{x}) + \lambda^{\text{coh}} E^{\text{coh}}(\mathbf{x}) \quad (2)$$

where $\lambda^{\text{app}}, \lambda^{\text{occl}}, \lambda^{\text{geom}}, \lambda^{\text{coh}}$ are scalar weights. We describe the energy terms below.

Function $E^{\text{app}}(\mathbf{x})$ favors correspondences between features having similar appearance. We define this function as a sum of unary terms:

$$E^{\text{app}}(\mathbf{x}) = \sum_{a \in A} \theta_a^{\text{app}} x_a. \quad (3)$$

For an assignment $a = (p', p'') \in A$, θ_a^{app} is the distance between appearance descriptors (such as Shape Context [14]) computed at points p' and p'' in the respective images. We have used different features depending on the task at hand (see sec. 4).

¹ The method in [23] first selects a *linear* function E^- which is an underestimator on the original objective function E , i.e. $E^-(\mathbf{x}) \leq E(\mathbf{x})$ for all feasible solutions \mathbf{x} . It then visits *all* feasible solutions \mathbf{x} with $E^-(\mathbf{x}) \leq E(\mathbf{x}^*)$ where $E(\mathbf{x}^*)$ is the cost of the optimal solution. For each solution a linear program is solved.

The term $E^{\text{occl}}(\mathbf{x})$ imposes a penalty for unmatched features. We define $E^{\text{occl}}(\mathbf{x})$ to be the fraction of unmatched features in the smallest of the two feature sets. We can write this function as

$$E^{\text{occl}}(\mathbf{x}) = 1 - \frac{1}{\min\{|P'|, |P''|\}} \sum_{a \in A} x_a \tag{4}$$

by noting that $\sum_{a \in A} x_a$ is equal to the number of distinct matched features in P' and P'' , $\forall \mathbf{x} \in M$. This result derives trivially from the uniqueness constraint.

The term $E^{\text{geom}}(\mathbf{x})$ is a measure of geometric compatibility between active correspondences. This term is similar to the distortion costs proposed in [17,18]. Note, however, that the energy terms used in these previous approaches include distortion costs for all pairs of matched features, which results in energy functions penalizing any deviation from a global rigid transformation. Instead, our function $E^{\text{geom}}(\mathbf{x})$ measures geometric compatibility of correspondences only for *neighboring* features. We demonstrate that this model permits more flexible mappings between the two sets of features and yields more accurate correspondences. We use a “neighborhood system” N to specify the pairs of correspondences involved in our measure of geometric compatibility. N consists of all correspondence pairs defined over neighboring features:

$$N = \{ \langle (p', p''), (q', q'') \rangle \in A \times A \mid p' \in N_{q'} \vee q' \in N_{p'} \vee p'' \in N_{q''} \vee q'' \in N_{p''} \} \tag{5}$$

where N_p indicates the set of K nearest neighbors of p (computed in the set of feature p), and K is a positive integer value controlling the size of the neighborhood, which we call *geometric neighborhood size*. $E^{\text{geom}}(\mathbf{x})$ is computed over pairs of active correspondences in the set N :

$$E^{\text{geom}}(\mathbf{x}) = \sum_{(a,b) \in N} \theta_{ab}^{\text{geom}} x_a x_b \tag{6}$$

where:

$$\theta_{ab}^{\text{geom}} = \eta(e^{\delta_{a,b}^2/\sigma_i^2} - 1) + (1 - \eta)(e^{\alpha_{a,b}^2/\sigma_\alpha^2} - 1) \tag{7}$$

$$\delta_{(p',p''),(q',q'')} = \frac{||p' - q'|| - ||p'' - q''||}{||p' - q'|| + ||p'' - q''||} \tag{8}$$

$$\alpha_{(p',p''),(q',q'')} = \arccos \left(\frac{p' - q'}{||p' - q'||} \cdot \frac{p'' - q''}{||p'' - q''||} \right) \tag{9}$$

Intuitively, $\theta_{(p',p''),(q',q'')}^{\text{geom}}$ computes the geometric agreement between neighboring correspondences $(p', p''), (q', q'')$ by evaluating how well the segment $\overline{p'q'}$ matches the segment $\overline{p''q''}$ in terms of both length and direction. The parameter η is a scalar value trading off the importance of preserving distances versus preserving directions.

The term $E^{\text{coh}}(\mathbf{x})$ favors spatial proximity of matched features. It incorporates our prior knowledge that matched features should form spatially coherent regions within each image, corresponding to common objects or parts in the image pair, in analogy to coherence on a pixel grid, used for example in image segmentation. We define the

cost $E^{\text{coh}}(\mathbf{x})$ as the fraction of neighboring feature pairs with different occlusion status (this can be viewed as an MRF Potts model over feature occlusion). We now show how to write this function directly in terms of solution \mathbf{x} . Let N_P be the set of pairs of neighboring features in the two images:

$$N_P = \{(p, q) \in (P' \times P') \cup (P'' \times P'') \mid p \in N_q \vee q \in N_p\}. \quad (10)$$

Then we can express $E^{\text{coh}}(\mathbf{x})$ as a sum of unary and pairwise terms:

$$E^{\text{coh}}(\mathbf{x}) = \frac{1}{|N_P|} \sum_{(p,q) \in N_P} V_{p,q}(\mathbf{x}) \quad (11)$$

where:

$$V_{p,q}(\mathbf{x}) = \sum_{a \in A(p)} x_a + \sum_{b \in A(q)} x_b - 2 \sum_{a \in A(p), b \in A(q)} x_a x_b. \quad (12)$$

$V_{p,q}(\mathbf{x})$ is equal to 0 if p, q are either both matched or both unmatched; $V_{p,q}(\mathbf{x})$ is equal to 1 otherwise.

Feature correspondence as graph matching. The problem defined above can be written as

$$\min_{\mathbf{x} \in M} E(\mathbf{x} \mid \bar{\theta}) = \sum_{a \in A} \bar{\theta}_a x_a + \sum_{(a,b) \in N} \bar{\theta}_{ab} x_a x_b \quad (13)$$

where the constraint set M is given by (1). This problem is often referred to as *graph matching* in the literature [22,5]. Features P' and P'' are viewed as vertices of the two graphs. Pairwise term $\bar{\theta}_{ab} x_a x_b$ with $a = (p', p'')$, $b = (q', q'')$ encodes compatibility between edges (p', q') , (p'', q'') of the first and second graph, respectively, while unary term $\bar{\theta}_a x_a$ measures similarity between vertices p', p'' .

We now address the question of how to optimize problem (13). Unfortunately, this problem is NP-hard [22]. We propose to use the *problem decomposition* approach (or *dual decomposition* - DD) for graph matching. Details are given in the next section.

3 Problem Decomposition Approach

On the high level, the idea is to decompose the original problem into several “easier” subproblems, for which we can compute efficiently a global minimum (or obtain a good lower bound). Combining the lower bounds for individual subproblems will then provide a lower bound for the original problem. The decomposition and the corresponding lower bound will depend on a parameter vector θ ; we will then try to find a vector θ that maximizes the bound. This approach is well-known in combinatorial optimization; sometimes it is referred to as “dual decomposition” [6]. It was applied to quadratic pseudo-boolean functions (i.e. functions of binary variables with unary and pairwise terms) by Chardaire and Sutter [7]. Their work is perhaps the closest to the method in this paper. As in [7], we use “small” subproblems for which the global minimum can be computed exactly in reasonable time. Our choice of subproblems for graph matching, however, is different from [7]. In vision the decomposition approach is probably

best known in the context of the MAP-MRF inference task. It was introduced by Wainwright et al. [8] who decomposed the problem into a convex combination of trees and proposed message passing techniques for optimizing vector θ . These techniques do not necessarily find the best lower bound. Schlesinger and Giginyak [9,10] and Komodakis et al. [11] proposed to use subgradient techniques [26,6] for MRF optimization, which guarantee to converge to a vector θ giving the best possible lower bound.

3.1 Graph Matching Via Problem Decomposition

We now apply this approach to the graph matching problem given by eq. (13). We decompose (13) into subproblems characterized by vectors θ^σ , $\sigma \in I$ with positive weights ρ_σ . (These weights are chosen a priori, and may affect the speed of convergence of the subgradient method in section 3.3.) Here I is a finite set of subproblem indexes. We will require the vector $\theta = (\theta^\sigma \mid \sigma \in I)$ to be a ρ -reparameterization of the original parameter vector $\bar{\theta}$ [8], i.e.

$$\sum_{\sigma \in I} \rho_\sigma \theta^\sigma = \bar{\theta} \tag{14}$$

For each subproblem σ we will define a lower bound $\Phi_\sigma(\theta^\sigma)$ which satisfies

$$\Phi_\sigma(\theta^\sigma) \leq \min_{\mathbf{x} \in M} E(\mathbf{x} \mid \theta^\sigma) \tag{15}$$

It is easy to see that the function

$$\Phi(\theta) = \sum_{\sigma \in I} \rho_\sigma \Phi_\sigma(\theta^\sigma) \tag{16}$$

is a lower bound on the original function. Indeed, if \mathbf{x}^* is an optimal solution of (13) then from (14)-(16) we get

$$\Phi(\theta) \leq \sum_{\sigma \in I} \rho_\sigma \min_{\mathbf{x} \in M} E(\mathbf{x} \mid \theta^\sigma) \leq \sum_{\sigma \in I} \rho_\sigma E(\mathbf{x}^* \mid \theta^\sigma) = E(\mathbf{x}^* \mid \bar{\theta})$$

In section 3.2 we will describe the subproblems that we use. For each subproblem σ we will do the following: (1) define constraints on vector θ^σ ; (2) define the function $\Phi_\sigma(\theta^\sigma)$; (3) specify an algorithm for computing $\Phi_\sigma(\theta^\sigma)$. In section 3.3 we will discuss how to maximize the lower bound $\Phi(\theta)$ using the subproblem solutions and, finally, how to obtain solution $\mathbf{x} \in M$ for our original problem.

3.2 Graph Matching Subproblems

Linear subproblem. In our first subproblem, which we denote by the index “ L ”, we require all pairwise terms to be zero: $\theta_{ab}^L = 0$ for $(a, b) \in N$. In such case problem (13) can be solved exactly in polynomial time, for example using the Hungarian algorithm [27]. (This is often known as the *linear assignment problem*.) We define $\Phi_L(\theta^L) = \min_{\mathbf{x} \in M} E(\mathbf{x} \mid \theta^L)$. To compute this minimum, we converted the problem

to an instance of a minimum cost circulation with unit capacities and ran the successive shortest path algorithm [27]. This solves the problem using $O(|P| + |A|)$ Dijkstra shortest path computations in graphs with $|P| + 1$ nodes and $O(|P| + |A|)$ edges.

Maxflow subproblem. In the second subproblem, which we denote by the index “ M ”, we do not put any restrictions on the vector θ^M . To get a lower bound, we ignore the uniqueness constraint $\sum_{a \in A(p)} x_a \leq 1$ and leave only the discreteness constraint: $x_a \in \{0, 1\}$. If the function $E(x | \theta^M)$ is submodular (i.e. coefficients θ_{ab}^M are non-positive for all pairwise terms $(a, b) \in N$), then we can compute a global minimum using a maxflow algorithm. With arbitrary θ_{ab}^M the problem becomes NP-hard [28]. We use the *roof duality* relaxation [29] to get a lower bound $\Phi_M(\theta^M)$ on the problem. It can be defined as the optimal value of the following linear program:

$$\Phi_M(\theta^M) = \min \sum_{a \in A} \theta_a^M x_a + \sum_{(a,b) \in N} \theta_{ab}^M x_{ab} \tag{17}$$

$$\text{subject to } \begin{cases} 0 \leq x_a \leq 1 & \forall a \in A \\ x_{ab} \leq x_a, \quad x_{ab} \leq x_b, \quad x_{ab} \geq x_a + x_b - 1, \quad x_{ab} \geq 0 & \forall (a, b) \in N \end{cases}$$

This relaxation can be solved in polynomial time by computing a maximum flow in a graph with $2(|A| + 1)$ nodes and $O(|A| + |N|)$ edges [28].

Local subproblems. For our last set of subproblems we use an exhaustive search to compute the global minimum (see [30] for details). Thus, we need to make sure that subproblems are sufficiently small. We use the following technique. For each point $p \in P$ we choose $N_p^d \subseteq P$ to be the set of K^d nearest points in the same image where K^d is a small constant, e.g. 2 or 3. (The superscript d stands for “decomposition”.) We then consider the subproblem which involves only assignments in the set $A(N_p^d) = \{(p', p'') \in A \mid p' \in N_p^d \vee p'' \in N_p^d\}$ and the edges between those assignments. More precisely, we require vector θ^p corresponding to this subproblem to satisfy the following constraints: (i) $\theta_a^p = 0$ if $a \notin A(N_p^d)$, and (ii) $\theta_{ab}^p = 0$ if $a \notin A(N_p^d)$ or $b \notin A(N_p^d)$. These constraints imply that we can fix assignments $a \in A - A(N_p^d)$ to 0 when computing the minimum $\min_{x \in M} E(x | \theta^p)$. Then we get a graph matching problem where the set of points in one of the images is N_p^d .

3.3 Algorithm Summary and Properties of Decomposition

Lower bound optimization. In the previous section we described constraints on vector θ and a lower bound $\Phi(\theta)$ consisting of $|P| + 2$ subproblems. It can be seen Φ is a concave function of θ . Furthermore, the constraints on θ yield a convex set Ω . This set is defined by the reparameterization equation (14) and constraints on individual subproblems $\theta^\sigma \in \Omega_\sigma$ given by equalities $\theta_i^\sigma = 0$, where index i here may denote either an assignment ($i = a$) or an edge ($i = (a, b)$). Let $I_i \subseteq I$ be the subsets of subproblem indexes for which θ_i^σ is **not** constrained to be 0. Thus, assignment $a \in A$ is involved in subproblems $\sigma \in I_a$, and edge $(a, b) \in N$ is involved in subproblems $\sigma \in I_{ab}$. Similar to [7,9,10,11], we used a projected subgradient method [26,6] for

maximizing $\Phi(\theta)$ over Ω . One iteration is given by $\theta := \mathcal{P}_\Omega(\theta + \lambda \mathbf{g})$ where \mathcal{P}_Ω is the operator that projects a vector to Ω , \mathbf{g} is a subgradient of $\Phi(\theta)$ and $\lambda > 0$ is a step size.

Projection. To project vector θ to Ω , we first compute vector $\hat{\theta} = \sum_\sigma \rho_\sigma \theta^\sigma$ and then update θ as follows: $\theta_i^\sigma := 0$ for $\sigma \in I - I_i$, and

$$\theta_i^\sigma := \theta_i^\sigma + \rho_\sigma \frac{\bar{\theta}_i - \hat{\theta}_i}{\sum_{\sigma' \in I_i} \rho_{\sigma'}^2} \quad \forall \sigma \in I_i$$

Subgradient computation. A subgradient of function $\Phi(\theta)$ is given by $\mathbf{g} = \sum_{\sigma \in I} \rho_\sigma \mathbf{g}^\sigma$ where \mathbf{g}^σ is a subgradient of function $\Phi_\sigma(\theta^\sigma)$. If the latter function is the global minimum of $E(\mathbf{x} | \theta^\sigma)$ (which is the case for $\sigma \in I - \{M\}$) then we can take $g_a^\sigma = x_a^\sigma$, $g_{ab}^\sigma = x_a^\sigma x_b^\sigma$ where \mathbf{x}^σ is a global minimizer of $E(\mathbf{x} | \theta^\sigma)$. For the maxflow subproblem a subgradient can be computed as $\mathbf{g}^M = \mathbf{x}^M$ where \mathbf{x}^M is an optimal solution of linear program (17). The method in [31] produces a half-integer optimal solution where $x_a^M \in \{0, 0.5, 1\}$ for all assignments a and x_{ab}^M is determined as follows: if $(x_a^M, x_b^M) \neq (0.5, 0.5)$ then $x_{ab}^M = x_a^M x_b^M$, otherwise $x_{ab}^M = 0$ if $\theta_{ab}^M \leq 0$ (i.e. the corresponding term is submodular) and $x_{ab}^M = 0.5$ if $\theta_{ab}^M > 0$.

Solution computation. To conclude the description of the method, we need to specify how to obtain solution $\mathbf{x} \in M$. We compute the solution in each iteration as follows: starting with labeling $\mathbf{x} = 0$, we go through local subproblems $\sigma \in I - \{L, M\}$ and assignments a involved in σ (in a fixed order), set $x_a = 1$ if $x_a^\sigma = 1$ and this operation preserves the uniqueness constraint on \mathbf{x} . (Here \mathbf{x}^σ denotes a global minimum of subproblem σ .) We maintain the solution with the smallest energy computed so far, and output it as a result of the method.

Further implementation details, e.g. the choice of the step size λ , are given in [30].

Properties of decomposition. It is not necessary to use all subproblems described in section 3.2. The only requirement is that each assignment $a \in A$ and edge $(a, b) \in N$ must be covered by at least one subproblem. In [30] we discuss how the choice of subproblems affects the optimal bound achievable with the decomposition method. Due to lack of space, here we give only a summary: (i) the linear subproblem is not essential with our choice of local subproblems; (ii) the maxflow subproblem is not essential if $K \leq K^d$; (iii) our lower bound is the same as or tighter than the bound of the roof duality approach [28] applied to an equivalent *quadratic pseudo-boolean optimization* formulation of problem (13).

4 Experimental Results

In most of our experiments we learned problem-specific parameters of our energy model from ground truth correspondences. We applied Nonlinear Inverse Optimization [20] (NIO) to learn non-negative parameters $\{\lambda^{\text{app}}, \lambda^{\text{occl}}, \lambda^{\text{geom}}, \lambda^{\text{coh}}, \eta, \sigma_l^2, \sigma_\alpha^2\}$. We used DD within NIO to optimize the learning objective (see [30] for details). The learning procedure was initialized using default parameters corresponding to uniform values for the weights $\{\lambda^{\text{app}}, \lambda^{\text{occl}}, \lambda^{\text{geom}}, \lambda^{\text{coh}}\}$, $\eta = 0.5$, and variance values $\sigma_l^2 = 0.5$, $\sigma_\alpha^2 = 0.9$.

In our experiments we compare the following algorithms:

DD. We used $K^d = \min\{K, 4\}$, where K is the geometric neighborhood size. Motivated by results in sec. 3.3, we did not use the linear subproblem. We set $\rho_\sigma = 1$ for all other subproblems σ . We used a maximum of 10000 iterations, and stopped earlier if the gap between the lower bound and the cost became smaller than 10^{-6} .

FUSION. This technique was introduced in [32] for MRF optimization with multiple labels. We propose to use it for graph matching as follows. First, we generate 256 solutions by applying one pass of coordinate descent (ICM) to zero labeling using random orders. (Different orders of visiting assignments usually yield different solutions.) We then “fuse” together pairs of solutions using the binary tree structure until a single solution remains. Fusion of solutions \mathbf{x}' , \mathbf{x}'' is defined as follows. First, we fix all assignments $a \in A$ for which \mathbf{x}' and \mathbf{x}'' agree, i.e. $x'_a = x''_a$. Then we convert the obtained graph matching problem to a quadratic pseudo-boolean optimization problem (see [30] for conversion details). Finally, we run the QPBO-PI method [33] starting either with labeling \mathbf{x}' if $E(\mathbf{x}' | \hat{\theta}) < E(\mathbf{x}'' | \hat{\theta})$ or with \mathbf{x}'' otherwise. The produced solution \mathbf{x} is guaranteed to have the same or smaller cost than the costs of \mathbf{x}' and \mathbf{x}'' .

BP. We converted graph matching to a quadratic pseudoboolean optimization problem and ran max-product belief propagation algorithm². We also tested applying the roof duality approach instead of BP, but results were quite discouraging (see details in [30]).

SMAC. We ran the spectral relaxation method of Cour et al. [12], using the graduated assignment algorithm [22] for discretization. Since SMAC imposes affine constraints on the solution, we applied this algorithm only to datasets without outliers, where the one-to-one affine constraint is satisfied. In principle, SMAC could handle outliers by the introduction of dummy nodes. However, this would increase the number of variables and potentially make the problem harder to solve.

COMPOSE. We reimplemented the algorithm in [13]. The problem was cast as assigning a label from the set $A(p') \cup \{\text{“occlusion”}\}$ to each point $p' \in P'$. Min-marginals for the linear subnetwork were computed via $O(|A| + |P'|)$ calls to the Dijkstra algorithm. As in [13], we used Residual Belief Propagation (RPB) [34] with damping=0.3 for computing pseudo min-marginals for the “smoothness” subnetwork containing pairwise terms $\theta_{ab}x_ax_b$. However, in our experiments messages did not converge, so we set an additional termination criterion for RPB: we stop it after passing $20|N|$ messages. As in [13], we computed the configuration by looking at individual messages at each node. We did not use damping for the outer loop since otherwise the produced configurations usually did not satisfy the uniqueness constraint.

HUNG. As in [5], we also tested the Hungarian algorithm using an energy consisting only of linear terms. On problems with occlusions, we used our occlusion cost in addition to the appearance energy term, i.e. $E^{\text{HUNG}}(\mathbf{x}) = \lambda^{\text{app}} E^{\text{app}}(\mathbf{x}) + \lambda^{\text{occl}} E^{\text{occl}}(\mathbf{x})$.

4.1 Comparative Results

Hotel sequence: wide baseline matching. We first demonstrate our approach on the CMU ‘hotel’ sequence³. We used as features the same manual labeling of 30 landmark

² We used the code from <http://www.adastral.ucl.ac.uk/~vladkolm/papers/TRW-S.html>

³ Available at: <http://vasc.ri.cmu.edu/idb/html/motion/hotel/index.html>.

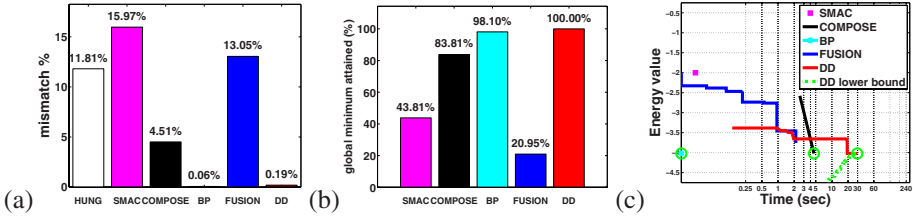


Fig. 1. Results on the Hotel sequence ($|P'| = |P''| = 30$, $|A| = 900$). (a) Mismatch percentages of HUNG and different optimizations applied to our energy model. (b) Frequency of convergence to global minimum. (c) Energy minimization versus time.

points employed in [5]. Since these points are visible in all frames and the motion is rigid, this matching problem is relatively simple. We include this dataset in our experiments as it was used in [5] to compare the performance of graph matching methods. We reproduce the experimental setup described in [5] using a subset of 105 frame pairs and adopting distances between Shape Context descriptors as unary terms. However, we replace the pairwise terms proposed in [5], with our geometric energy function $E^{\text{geom}}(x)$, using $K = 2$. Due to the absence of outliers, we remove $E^{\text{coh}}(x)$ from our energy and use a large constant value for λ^{occl} . We set the remaining parameters to default values, as defined above. We set $A = P' \times P''$. Figure 1(a) shows the matching error obtained by optimizing this model with different methods. We include in the plot also the performance of HUNG. Here BP and DD perform dramatically better than the other methods, with errors approaching 0%. Note that the error of our system is over 50 times smaller than the errors reported in [5]. On this dataset DD found always the global minimum within a minute (see Figure 1(b)). Figure 1(c) illustrates performance versus runtime on one image pair (frame 1 and 64). In this plot we indicate convergence to a global minimum with a green circle. BP does well on this sequence, nearly matching the minimization performance of DD, at a reduced cost. We also implemented the energy function described in [5]. This model uses Delaunay triangulation to define the graph topology in each image, and employs binary edge compatibility values in $\{0,1\}$. DD provided the best performance with this model, with a 5.7% error. This suggests that both our model and our optimization contribute to the improvement over the results reported in [5] where the best system had a matching error above 10%.

Matching MNIST digits. Here we describe experiments on images of handwritten digits from the MNIST dataset [35]. For training, we randomly sampled from this dataset one image pair for each of the 10 digits. We repeated the same procedure to generate a test set of 10 pairs of same digits. From each pair we extracted point sets P' and P'' by uniformly sampling 100 points along the Canny edges of each image, using the procedure described in [14]. We defined the unary potentials $\theta_{(p',p'')}^{\text{app}}$ to be the Euclidean distances between Shape Context descriptors computed at points p' , p'' . We formed the set of candidate assignments $A \in P' \times P''$ by selecting the 5 most similar features, in terms of Shape Context distance, for each point $p \in P$. We collected ground truth correspondences in the set $(P' \times P'')$ for each of the 20 image pairs. The parameters of our model were learned from the 10 training image pairs with NIO. Figure 2(a) shows that

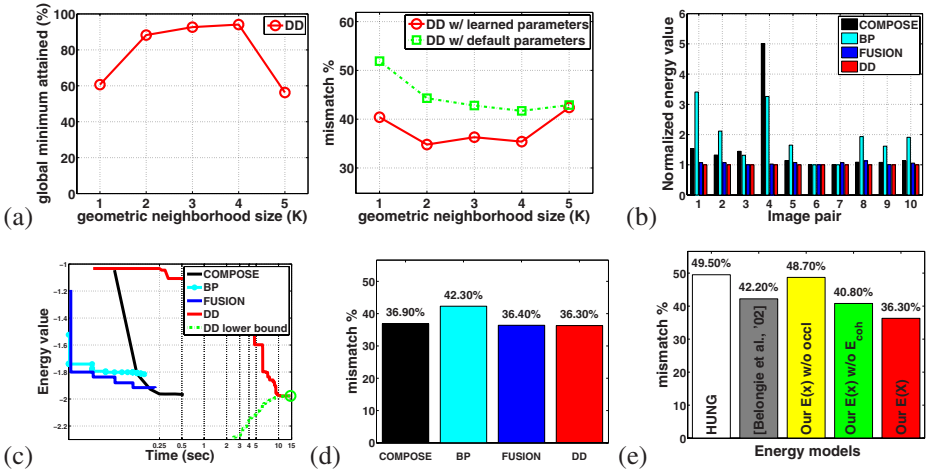


Fig. 2. Experimental results on MNIST digits ($|P'| = |P''| = 100$, $|A| = 695$, on average). (a) Correlation between learning accuracy and matching performance: the left plot shows the frequency of global minimum convergence during learning versus K ; the right plot shows mismatch error on test set. (b) Normalized energy values. (c) Optimization performance versus runtime. (d) Mismatch error comparison between different optimization methods using our energy model. (e) Mismatch error using different energy models.

the matching accuracy on the test set critically depends on the ability to globally optimize the energies during the model learning stage. The left plot reports the frequency of convergence to a global minimum during learning, plotted as a function of K , the geometric neighborhood size. The second plot shows the test set matching error of DD with learned versus default parameters. Matching error here is measured as percentage of incorrect correspondences (as defined in [30]). We can see that the matching is much more accurate when using the parameters for which DD reached more frequently global optimality during learning. Interestingly, although the frequency of global minimum convergence increases slightly when varying K from 2 to 4, the matching error remains roughly the same. This suggests that geometric penalty terms defined over small neighborhoods are sufficient to spatially regularize the correspondences. Thus, models involving geometric costs defined over all pairs of matched features, such as those used in [17,18], may be unnecessarily restrictive for many applications, in addition to being more difficult to optimize.

Given these results, we have used the model learned with $K = 3$ for the MNIST experiments described below. Figure 2(b) shows the normalized energy values obtained by different optimization methods on the test set. For each family of results we performed an *additive* normalization so that for each image pair the energy of the best method becomes a fixed number. On 9 out of the 10 test image pairs, DD reaches global optimality, and provides the minimum energy value on all examples. FUSION, BP, and COMPOSE find the global minimum only on 2 images. FUSION finds solutions with energy values very close to those obtained by DD. COMPOSE and BP provide considerably higher

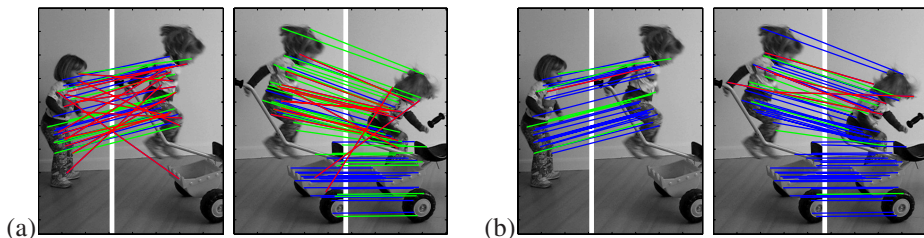


Fig. 3. Estimating human motion ($|P'| = 118$, $|P''| = 172$, $|A| = 1128$ on average). Correspondences computed with (a) the Hungarian method and (b) DD. Correct correspondences are shown in blue, missed assignments in green, and mismatches in red.

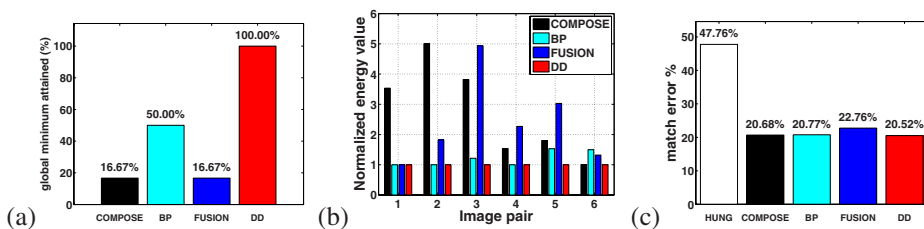


Fig. 4. Experiments on human motion frames. (a) Frequency of convergence to global minimum. (b) Normalized energy values. (c) Correspondence error.

energy values on some of the examples. Figure 2(c) shows minimization performance as a function of time, evaluated on a sample image pair. Figure 2(d) shows the correspondence accuracy obtained by optimizing our energy with the different methods. Again, we find that DD and FUSION yield the best accuracy. We also evaluated variations of the energy model defined in Equation (2) obtained by using only the linear terms (HUNG), by dropping the spatial coherence term, and by forcing all points to be matched (implemented by fixing λ^{occl} to a large value). The parameters of these modified models were learned again with NIO, using DD for both training and testing. We see from Figure 2(e) that both the spatial coherence prior, as well as the occlusion cost, improve the matching accuracy. On these instances the simple appearance-based model used by HUNG gives poor accuracy. We also report the matching error given by the model and optimization method of Belongie et al., which was applied to MNIST digit examples in [14]. Our approach performs better than this state-of-the-art method.

Estimating long range non-rigid motion. In this subsection we describe results on the task of estimating large-disparity motion. For this experiment we used four (time-separated) video frames of a child jumping. We matched each image to every other image, for a total of six matches. The motion between any pair of these pictures is very large and highly non-rigid. There is self-occlusion created by the motion of arms and torso, and occlusion due to a tricycle positioned between the child and the camera. Feature points were extracted by running the Harris corner detector on each image. We used

Euclidean distances of geometric blur descriptors [17] computed at each feature point, both for selecting assignments in A (by choosing the five most similar features for each point $p \in P$) as well as for calculating the unary terms of our energy. We learned the parameters in our model by applying the NIO algorithm to ground truth correspondences of two image pairs from a separate sequence containing the same child walking. Here we report results using $K = 6$. Figure 3 shows two matching examples from this experiment and correspondences found with HUNG and DD. Note the ability of our system to cope well with occlusion and multiple motions. DD converged to a global minimum on all the image pairs in this experiment (see Figure 4(a,b)). Figure 4(c) reports the correspondence errors (including mismatches as well as missed assignments).

Additional results and experiments on another dataset are given in [30].

5 Conclusions

We have introduced novel models and optimization algorithms for feature correspondence. We believe to be the first to demonstrate graph matching techniques capable of reaching global optimality on various real-world image matching problems. As a future work, we plan to replace exhaustive search for local subproblems with a branch-and-bound method, as in [7]. We hope that this may speed up substantially the DD method.

Acknowledgements. We are grateful to Timothee Cour, Praveen Srinivasan, and Jianbo Shi for providing the software of their SMAC algorithm. We thank John Duchi, Gal Elidan and Danny Tarlow for sharing code and answering questions about the COMPOSE method. Thanks to Tiberio Caetano for providing the Hotel feature data.

References

1. Belhumeur, P.N.: A binocular stereo algorithm for reconstructing sloping, creased, and broken surfaces in the presence of half-occlusion. In: ICCV (May 1993)
2. Kolmogorov, V., Zabih, R.: Computing visual correspondence with occlusions using graph cuts. In: ICCV (2001)
3. Dorko, G., Schmid, C.: Selection of scale-invariant parts for object class recognition. In: ICCV (2003)
4. Sivic, J., Russell, B., Efros, A., Zisserman, A.: Discovering object categories in image collections. In: ICCV (2005)
5. Caetano, T.S., Cheng, L., Le, Q.V., Smola, A.J.: Learning graph matching. In: ICCV (2007)
6. Bertsekas, D.: Nonlinear Programming. Athena Scientific (1999)
7. Chardaire, P., Sutter, A.: A decomposition method for quadratic zero-one programming. *Management Science* 41(4), 704–712 (1995)
8. Wainwright, M.J., Jaakkola, T.S., Willsky, A.S.: MAP estimation via agreement on trees: Message-passing and linear-programming approaches. *IEEE Trans. Information Theory* 51(11) (2005)
9. Schlesinger, M.I., Giginyak, V.V.: Solution to structural recognition (MAX,+)-problems by their equivalent transformations. Part 1. *Control Systems and Computers* (1), 3–15 (2007)
10. Schlesinger, M.I., Giginyak, V.V.: Solution to structural recognition (MAX,+)-problems by their equivalent transformations. Part 2. *Control Systems and Computers* (2), 3–18 (2007)

11. Komodakis, N., Paragios, N., Tziritas, G.: MRF optimization via dual decomposition: Message-passing revisited. In: ICCV (2007)
12. Cour, T., Srinivasan, P., Shi, J.: Balanced graph matching. In: NIPS (2007)
13. Duchi, J., Tarlow, D., Elidan, G., Koller, D.: Using combinatorial optimization within max-product belief propagation. In: NIPS (2007)
14. Belongie, S.J., Malik, J., Puzicha, J.: Shape matching and object recognition using shape contexts. *PAMI* 24(4), 509–522 (2002)
15. Torr, P.H.S.: Geometric motion segmentation and model selection. *Philosophical Transactions of the Royal Society*, 1321–1340 (1998)
16. Sclaroff, S., Pentland, A.: Modal matching for correspondence and recognition. *PAMI* 17(6), 545–561 (1997)
17. Berg, A., Berg, T., Malik, J.: Shape matching and object recognition using low distortion correspondence. In: CVPR (2005)
18. Leordeanu, M., Hebert, M.: A spectral technique for correspondence problems using pairwise constraints. In: ICCV (2005)
19. Torr, P.H.S.: Solving Markov random fields using semi definite programming. In: AISTATS (2003)
20. Liu, C.K., Hertzmann, A., Popović, Z.: Learning physics-based motion style with nonlinear inverse optimization. *ACM Trans. on Gr.* 24(3), 1071–1081 (2005)
21. Conte, D., Foggia, P., Sansone, C., Vento, M.: Thirty years of graph matching in pattern recognition. *Int. J. of Pattern Recognition and Artificial Intelligence* 18(3), 265–298
22. Gold, S., Rangarajan, A.: A graduated assignment algorithm for graph matching. *PAMI* 18(4), 377–388 (1996)
23. Maciel, J., Costeira, J.: A global solution to sparse correspondence problems. *PAMI* 25(2), 187–199 (2002)
24. Cabot, A., Francis, R.: Solving certain nonconvex quadratic minimization problems by ranking the extreme points. *Operations Research* 18(1), 82–86 (1970)
25. Schellewald, C., Schnörr, C.: Probabilistic subgraph matching based on convex relaxation. In: EMMCVPR (2005)
26. Shor, N.Z.: *Minimization methods for nondifferentiable functions*. Springer, Heidelberg (1985)
27. Ahuja, R., Magnanti, T., Orlin, J.: *Network Flows: Theory, Algorithms, and Applications*. Prentice Hall, Englewood Cliffs (1993)
28. Boros, E., Hammer, P.: Pseudo-boolean optimization. *Discr. Appl. Math.* 123(1-3) (2002)
29. Hammer, P.L., Hansen, P., Simeone, B.: Roof duality, complementation and persistency in quadratic 0-1 optimization. *Mathematical Programming* 28, 121–155 (1984)
30. Torresani, L., Kolmogorov, V., Rother, C.: Feature correspondence via graph matching: Models and global optimization. Technical Report MSR-TR-2008-101 (2008)
31. Boros, E., Hammer, P.L., Sun, X.: Network flows and minimization of quadratic pseudo-Boolean functions. Technical Report RRR 17-1991, RUTCOR (May 1991)
32. Lempitsky, V., Rother, C., Blake, A.: LogCut - efficient graph cut optimization for Markov random fields. In: ICCV (2007)
33. Rother, C., Kolmogorov, V., Lempitsky, V., Szummer, M.: Optimizing binary MRFs via extended roof duality. In: CVPR (2007)
34. Elidan, G., McGraw, I., Koller, D.: Residual belief propagation: Informed scheduling for asynchronous message passing. In: UAI (2006)
35. LeCun, Y., Bottou, L., Bengio, Y., Haffner, P.: Gradient-based learning applied to document recognition. *Proceedings of the IEEE* 86(11), 2278–2324 (1998)