# Probabilistic Region Matching in Narrow-Band Endoscopy for Targeted Optical Biopsy

Selen Atasoy[1,2], Ben Glocker[1], Stamatia Giannarou[2], Diana Mateus[1],
Alexander Meining(MD)[3], Guang-Zhong Yang[2] and Nassir Navab[1],

[1] Computer Aided Medical Procedures (CAMP), Technische Universität München
{atasoy, glocker, mateus}@cs.tum.edu
[2] Visual Information Processing Group, Imperial College London
{catasoy, stamatia.giannarou, g.z.yang}@imperial.ac.uk
[3] Department of Gastroenterology, Technische Universität München
{alexander.meining}@lrz.tu-muenchen.de

**Abstract.** Recent advances in biophotonics have enabled *in-vivo, in-situ* histopathology for routine clinical applications. The non-invasive nature of these optical 'biopsy' techniques, however, entails the difficulty of identifying previously visited biopsy locations, particularly for surveillance examinations. This paper presents a novel region-matching approach for narrow-band endoscopy to facilitate retargeting the optical biopsy sites. The task of matching sparse affine covariant image regions is modelled in a Markov Random Field (MRF) framework. The proposed model incorporates appearance based region similarities as well as spatial correlations of neighbouring regions. In particular, a geometric constraint that is robust to deviations in relative positioning of the detected regions is introduced. In the proposed model, the appearance and geometric constraints are evaluated in the same space (photometry), allowing for their seamless integration into the MRF objective function. The performance of the method as compared to the existing state-of-the-art is evaluated on both *in-vivo* and simulation datasets with varying levels of visual complexities.
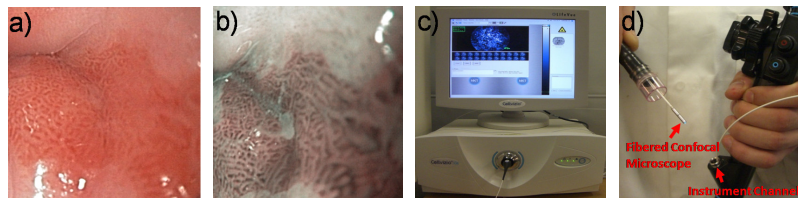
**Keywords:** In-vivo histology, fibered confocal microscopy, narrow-band endoscopy, Markov Random Fields, affine covariant regions.

## 1 Introduction

Oesophageal Adenocarcinoma (OAC) is the most rapidly increasing cancer in Europe and the United States, which has a 5-year survival rate of only 10% [1]. Barrett's Oesophagus (BO), referring to the abnormal change of the oesophageal mucosa caused by gastro-oesophageal reflux (Fig. 1a-b), is the only recognized precursor to OAC. Therefore, for patients diagnosed with BO, periodic surveillance by gastrointestinal (GI) endoscopy together with systematic biopsy is important for the early detection and prevention of OAC.

In current surveillance protocols, a new technique called Narrow-Band Endoscopic Imaging (NBI), has shown advantage compared to conventional white light

endoscopy as it allows for detailed visualization of mucosa and the underlying vascular patterns (Fig. 1a-b). A further technique called Fibered Confocal Microscopy (FCM) which enables real-time visualization of cellular structures *in-vivo* and *in-situ* (Fig. 1c) is also introduced recently. During GI endoscopy, a fibered confocal microprobe can be inserted easily through the instrument channel of a standard endoscope (Fig. 1d), providing *in-situ* histopathology without the need for tissue biopsy [2]. This has significant benefits in terms of ease of examination, patient comfort and real-time feedback. In practice, however, the non-invasive nature of the procedure also makes it difficult to return to previously examined biopsy sites in surveillance endoscopy due to the absence of scar on the tissue. The purpose of this paper is to present a novel image-based region matching method for biopsy site re-targeting in NBI.



**Fig. 1.** Appearance of BO a) in white light endoscopy and b) in NBI. c) The FCM machine and d) FCM microprobe passing through the instrument channel of a standard endoscope.

Region matching in NBI entails several challenges, which include tissue deformation, prevalence of similar surface textures and mucosal patterns. As the endoscope is very close to the tissue, small differences in the visible scales of the same feature can cause a significant change in the visual content. Furthermore, the common issue of view-invariant scene matching also needs to be addressed.

Viewpoint invariant scene matching is a well studied problem in computer vision and it typically proceeds by representing the scene as a collection of affine covariant regions which are described by a vector computed from the regions' appearances. Usually a nearest neighbour matching of the descriptor vectors incorporating geometric constraints is used to eliminate possible outliers ([3-5] and references within). In endoscopy, the major focus is directed towards short-baseline matching/tracking in the presence of tissue deformation [6].

Recently, spectral methods have been proposed for region-matching in images undergoing non-rigid transformations [7-9]. These methods model a graph for the feature set in each image and estimate their correspondences by graph matching. Thereby the geometric relations are modelled in terms of point locations where distance and/or orientation preservation is enforced. The main focus of these approaches lies on estimating the optimal solution for the NP-hard graph matching problem rather than on optimal modelling. However, Caetano *et al.* have demonstrated that finding the (near) optimal graphs can greatly simplify the matching problem and improve the results [10]. The authors proposed a learning based approach for optimal graph extraction. An MRF model [11] and a graph matching approach [8] with optimal model parameter learning are also presented for the correspondence problem.
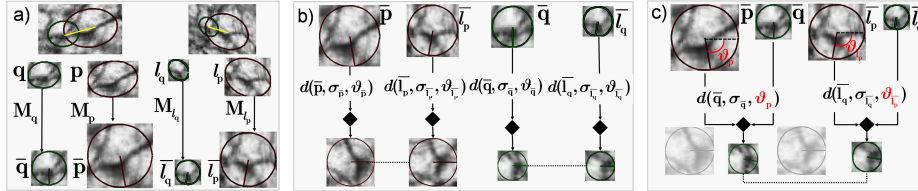
In this paper, we will focus on deriving the (near) optimal MRF model for the feature correspondence problem in NBI. The proposed model incorporates appearance based region similarities as well as the spatial correlations of neighbouring regions. To this end, we introduce a geometric constraint that evaluates the consistency between neighbouring matches on their photometric properties. Evaluation of the appearance and geometric constraints in the same space (photometry) allows for their seamless integration. The performance of the proposed method is evaluated with both *in-vivo* and simulation datasets.

## 2 Methods

The proposed method involves initial affine covariant region detection. This is followed by formulating a particular MRF model for the matching problem. Finally, the optimal labelling is computed using Belief Propagation.

### 2.1 Region Detection and Description

Affine covariant regions are detected independently on both images using affine invariant anisotropic region detector [12], which is shown to be robust against small deformations. For viewpoint invariant description, each elliptical region $p$ is normalized by the corresponding affine transformation $M_p$ (determined by the shape of the ellipse) and mapped onto the corresponding circular region $\bar{p} = M_p \cdot p$ (Fig. 2a). Then, the dominant gradient orientation $\vartheta_{\bar{p}}$ is estimated from the local image gradients and the SIFT descriptor [3] $d(\bar{p}, \sigma_{\bar{p}}, \vartheta_{\bar{p}})$ is computed from the circular patch $\bar{p}$ using the characteristic scale $\sigma_{\bar{p}}$ and the dominant gradient orientation $\vartheta_{\bar{p}}$.



**Fig. 2.** a) Viewpoint invariant region description. b) Unary costs computed from the region descriptors, where the diamond indicates the SIFT descriptor computed on the affine normalized patches. c) The proposed pair-wise costs. SIFT descriptors are computed on the patches $\bar{q}$ and $\bar{l}_q$ using the dominant gradient orientations $\vartheta_p$, $\vartheta_{l_p}$ of the regions $\bar{p}$ and $\bar{l}_p$. Two compared image patches are the same, whereas the length and orientation of the line segment between two region centres are not preserved as illustrated via the yellow lines in a).

### 2.2 Matching through Markov Random Fields

Given the computed region descriptors, we model the matching problem as global optimization of an MRF labelling. We define the regions in the first image to be the nodes $G = \{1,..,n\}$ of the MRF and the regions in the second image to be the labels $\mathcal{L}^+ = \{l_0, l_1,.., l_M\}$ including the null-label $l_0$, which is assigned to regions without true correspondence in the second image. In this paper, we consider only up to pair-

wise relations. Thus, finding the maximum a posteriori (MAP) estimate of the optimum labelling $\mathbf{l}^*$ is equivalent to minimizing the energy function:

$$E_{\mathrm{MRF}}(\mathbf{l}) = \sum_{p \in \mathcal{G}} V_p(l_p) + \sum_{p \in \mathcal{G}} \sum_{q \in \mathcal{N}(p)} V_{pq}(l_p, l_q) \tag{1}$$

where $V_p(l_p)$ is the unary cost of assigning the label $l_p$ to the node $p$, $V_{pq}(l_p, l_q)$ is the pair-wise cost and $\mathcal{N}$ defines the neighbourhood system.

## 2.3    Unary Costs

In our model, photometric similarities between the node and the label regions are evaluated via the unary costs by defining $V_p(l_p)$ to be the distance of the SIFT descriptors of the node $\overline{p}$ and label $\overline{l}_p$ regions (Fig. 1b). We further define the cost $V_p(l_0)$ of assigning the null-label $l_0$ to a node $p$ to be a function of the photometric similarities. The motivation is that assigning the null-label $l_0$ to a region that has a strong correspondence in the second image should have a higher cost than assigning it to a region with no (strong) correspondence. We define the *null-cost function* of the node $p$ as $V_p(l_0) = \alpha \cdot (1 - \min(V_p(\cdot)))$, where $\min(V_p(\cdot))$ is the minimum cost of assigning a label to the node $p$, and $\alpha$ is the factor regulating the trade-off between the quality and the number of matches. (For all our *in-vivo* datasets, the best performance is achieved for $\alpha = 0.5$). The final unary costs are computed as:

$$V_p(l_p) = \begin{cases} \arccos(d(\overline{p}, \sigma_{\overline{p}}, \vartheta_{\overline{p}}) \cdot d(\overline{l}_p, \sigma_{\overline{l}_p}, \vartheta_{\overline{l}_p})) / \arccos(0) & \text{if } l_p \neq l_0 \\ \alpha \cdot (1 - \min(V_p(\cdot))) & \text{otherwise} \end{cases}, \tag{2}$$

where all costs $V_p(\cdot)$ are normalized to the interval $[0,1]$ by dividing by the maximum possible angle between two descriptor vectors; $\arccos(0)$.

## 2.4    Neighbourhood Systems

In the context of the matching problem, each region is allowed to have at most one correspondence in the second image, *i.e.*, each label can be assigned at most to one node. This uniqueness constraint is included into the energy function by connecting each node with all the other nodes within the *global neighbourhood system* $\mathcal{N}$ and by defining the pair-wise cost for assigning the same label to two nodes to be infinite: $V_{pq}(l_p, l_q) = \infty$ if $l_p = l_q \neq l_0$. We further define a *local neighbourhood system* in order to impose flexible local geometric constraints. This neighbourhood system is defined for both the nodes and the labels to impose neighbourhood preservation as the initial geometric constraint (Eq. 3.1). The local neighbourhood $\mathcal{N}_{\mathrm{local}}(p)$ of a region $p$ is set to be: $\mathcal{N}_{\mathrm{local}}(p) = \{q \neq p \mid \; \|p - q\| < t\}$, where $\|p - q\|$ is the Euclidean distance between the centres of $p$ and $q$ and $t$ is a threshold value. (We use $t = 10\%$ and $t = 20\%$ of the image size for the node- and label-neighbourhoods

respectively to ensure the connectivity of two neighbouring regions after a large viewpoint change). For regions fulfilling the neighbourhood preservation, a novel geometric constraint is imposed measuring the consistency of two matches.

## 2.5 Pair-wise Costs

In this paper, we propose a geometric constraint based on the assumption that neighbouring regions move with similar transformations. The idea is as follows: if two neighbouring regions $p$ and $q$ have the corresponding regions $l_p$ and $l_q$ in the second image, then there exist two affine transformations $A_p$ and $A_q$ such;
$$l_p(\mathbf{x}) = A_p \cdot p(\mathbf{x}) = s_p \cdot R_p \cdot M_p \cdot p(\mathbf{x}) \quad \text{and} \quad l_q(\mathbf{x}) = A_q \cdot q(\mathbf{x}) = s_q \cdot R_q \cdot M_q \cdot q(\mathbf{x})$$
where $s_p$ and $s_q$ are scale factors. Theoretically, for spatially close regions on the same plane it holds $A_p = A_q$. However, for neighbouring regions on different planes this assumption is too restrictive and can be relaxed by assuming only $R_p = R_q = R$, where $R$ is the rotation of the local neighbourhood between two images.

If two neighbouring matches $m_p = (p, l_p)$ and $m_q = (q, l_q)$ are true correspondences, then the SIFT descriptors $d(\overline{q}, \sigma_{\overline{q}}, \vartheta_{\overline{p}})$ and $d(\overline{l_q}, \sigma_{\overline{l_q}}, \vartheta_{\overline{l_p}})$ computed on the patches $\overline{q}$ and $\overline{l_q}$ using their own characteristic scales $\sigma_{\overline{q}}$, $\sigma_{\overline{l_q}}$ and the dominant gradient orientations $\vartheta_{\overline{p}}$, $\vartheta_{\overline{l_p}}$ of $\overline{p}$ and $\overline{l_p}$ should be similar as the local rotation can be determined as $R = \vartheta_{\overline{l_p}} - \vartheta_{\overline{p}}$ (Fig. 1c). (Recall that $M_p \cdot p = \overline{p}$ and $M_q \cdot q = \overline{q}$). This similarity measure indicates the consistency of the matches $m_p$ and $m_q$, as the rotation estimated from the match $m_p$ is evaluated on the regions of the match $m_q$. Combining with the neighbourhood preservation, the pair-wise costs to evaluate geometric constraints are defined as:

$$\psi_{pq}(l_p, l_q) = \begin{cases} \infty & \text{if } (l_p \notin \mathcal{N}_{\text{local}}(l_q)) \quad \textbf{(3.1)} \\ \arccos(d(\overline{q}, \sigma_{\overline{q}}, \vartheta_{\overline{p}}) \cdot d(\overline{l_q}, \sigma_{\overline{l_q}}, \vartheta_{\overline{l_p}})) / \arccos(0) & \text{if } (l_p \in \mathcal{N}_{\text{local}}(l_q)) \quad \textbf{(3.2)} \end{cases}$$

Introducing the geometric constraints the final pair-wise costs are defined as:

$$V_{pq}(l_p, l_q) = \begin{cases} \infty & \text{if } (l_p = l_q \neq l_0) \\ \alpha & \text{if } (l_p = l_0) \vee (l_q = l_0) \\ \psi_{pq}(l_p, l_q) & \text{if } (q \in \mathcal{N}_{\text{local}}(p)) \\ 0 & \text{otherwise} \end{cases} \quad \textbf{(4)}$$

where $\alpha$ is the same factor used in the unary null-costs.

The derived geometric constraint (Eq. 3.2) is invariant to changes in the scale and relies only on the assumption of similar rotations within a local neighbourhood. This is in contrast to previous methods taking global consistency of the features into account. Such methods would fail in case of global deformation, which is present in our applications. This constraint also allows us to be locally more flexible than recent methods making stronger assumptions, such as invariance of the distance between the

neighbouring points or the orientation of line segment connecting their centroids [7-9]. Furthermore, the proposed constraint evaluates local image geometry based on the photometric properties of the patches rather than their spatial locations. This allows for the evaluation of the unary and the pair-wise costs in the same space and combining them in the objective function without using any weighting parameters.
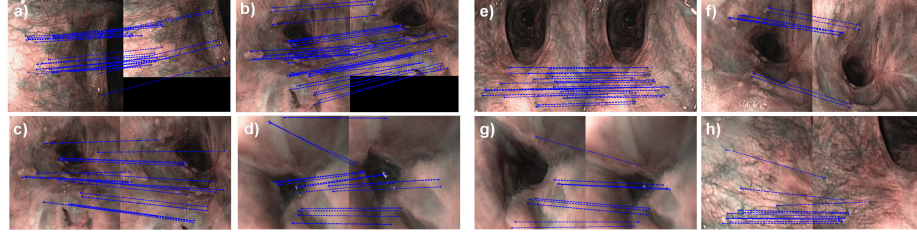
## 2.6 MAP Estimation

The MAP labelling of the proposed MRF model is estimated using Belief Propagation (BP) [13]. The non-submodularity of the pair-wise costs restricts the choice of the MRF inference algorithms to those without prior constraints on the class of energy functions. In this paper, without loss of generality we use the BP algorithm.
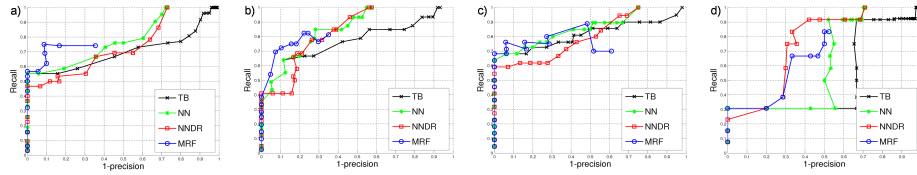
## 3 Experiments and Results

The performance of the proposed method is evaluated on 4 *in-vivo* and 4 simulation datasets and compared to 3 matching strategies evaluated in [14]. The regions are detected and described as explained in Section 3.1. The threshold-based (TB) [14], nearest-neighbour (NN) [14] and the nearest neighbour distance ratio matching (NNDR) [3,14] are applied for varying threshold values. MRF-based method is performed for different values of the factor $\alpha$ within the convergence range of the optimization. We further compared the hypergraph matching algorithm (HGM) using the proposed affine invariant geometric measure via quadripartite point relations [7]. However, these graph matching methods are not adapted to large number of non-matching regions (43%-87% in our datasets). Therefore, the performance of the HGM was poor and is not illustrated here in detail. For quantitative analysis, we evaluate $recall$ (the ratio of correct matches to the total number of correspondences) versus $1 - precision$ (the number of false matches with respect to the number of matched regions). For the best matching results $recall = 1$ and $1 - precision = 0$.

### 3.1 Simulation Studies

For evaluation with known ground truth data, we created 4 simulation datasets (2 for viewpoint change and 2 for deformation). In the first study, we generated images under different viewpoint conditions by transforming 2 *in-vivo* images (one veined and one structured area) with known homographies. In the second study, we deformed 2 *in-vivo* images (one structured and one homogenous area) and tracked the detected regions. In both studies, two regions were accepted as a correct match if the distance between the centres of the transformed and detected ellipses was less than 1% of the image size and the overlap was more than 55%. Figs. 4(a-c) demonstrate that MRF matching results in a better performance than all compared methods for structured scenes. Fig. 4d shows that in the presence of non-distinctive regions, MRF-matching and NNDR (which favours distinctive matches) exhibit a similar performance. The matching results of the MRF model are presented in Figs. 3(a-d).
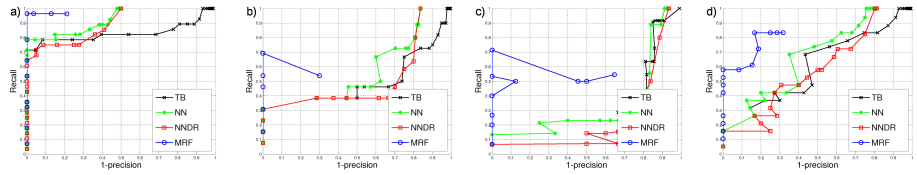
**Fig. 3.** Matching results of the MRF model on a-b) viewpoint change c-d) deformation simulation datasets e) first, f) second, g) third and h) fourth *in-vivo* datasets.



**Fig. 4.** Validation of the results on simulation datasets. Viewpoint change on images of a) veined b) structured tissue. Deformation on images of c) structured d) homogenous tissue.

## 3.2 In-vivo Studies

For the *in-vivo* studies, we used 4 NBI datasets with different viewpoint and photometric conditions. The first 3 datasets contain two distant frames of the same GI procedure from different viewpoints showing a veined area (Fig. 3e), structured area (Fig. 3f) and homogenous area with large deformation (Fig. 3g). The fourth dataset contains images acquired during two different GI examinations with a time difference of 3 months where the patient underwent chemotherapy. This results in large changes in the visual appearance of the tissue (Fig. 3h). For the *in-vivo* data sets, the ground truth data was provided by manual labelling. Fig. 5 demonstrates that for all *in-vivo* cases the proposed MRF model performs better than the state-of-the art descriptor matching techniques. For all datasets (simulation and *in-vivo*) maximum recall values for the acceptable precision interval (80%-100% inliers) are summarized in Table1. The matching results for the *in-vivo* datasets are presented in Fig 3e-f.



**Fig. 5.** Validation of the results on *in-vivo* datasets. a-b-c-d) show the recall versus precision of each matching algorithm for the first, second, third and fourth in-vivo dataset, respectively.

## 4 Conclusion

In this paper, we investigate the task of region matching for NBI and propose a new method towards an image-based solution for consistent re-targeting of optical biopsy sites. To this end, we present an MRF model for matching affine covariant regions

incorporating a novel geometric constraint for dealing with large changes in the observed datasets. Our results demonstrate the robustness of the proposed model for deformable wide-baseline matching on *in-vivo* and simulation datasets. For future work, we plan to further extend our approach in order to provide a complete framework for this novel and challenging application. This would require both detection and tracking taking into account the sequential appearances of features within the first and secondary examinations.

|  | *Viewp.1* | *Viewp.2* | *Def.1* | *Def.2* | *In-vivo1* | *In-vivo2* | *In-vivo3* | *In-vivo4* |
|---|---|---|---|---|---|---|---|---|
| **MRF** | 0.75 | 0.75 | 0.76 | 0.31 | 0.96 | 0.69 | 0.71 | 0.83 |
| **NN** | 0.58 | 0.66 | 0.73 | 0.31 | 0.82 | 0.31 | 0.13 | 0.37 |
| **NNDR** | 0.53 | 0.68 | 0.62 | 0.31 | 0.75 | 0.31 | 0.06 | 0.21 |
| **TB** | 0.55 | 0.63 | 0.73 | 0.31 | 0.78 | 0.31 | 0.13 | 0.32 |

**Table 1.** Summary of the maximum recall values for the precision interval [0.8-1.0] (80%-100% inliers) for the simulation and *in-vivo* datasets.

# References

1. Wani, S. and Sharma, P.: The rationale for screening and surveillance of Barrett's metaplasia, Best Practice & Research Clinical Gastroenterology, 20, 829–842 (2006)
2. Meining, A., Saur, D., Bajbouj, M., Becker, V., Peltier, E., Höfler, H., von Weyhern, C., Schmid, R. and Prinz, C.: In Vivo Histopathology for Detection of Gastrointestinal Neoplasia With a Portable, Confocal Miniprobe: An Examiner Blinded Analysis Clinical Gastroenterology and Hepatology, Elsevier, 5, 1261–1267, (2007)
3. Lowe, D.: Distinctive Image Features from Scale-Invariant Keypoints, Int. Journal of Computer Vision, 91–110, (2004).
4. Sivic, J. and Zisserman, A.: Video Google: a text retrieval approach to object matching in videos, ICCV, 1470–1477, (2003)
5. Schaffalitzky, F. and Zisserman, A. Automated location matching in movies, Computer Vision and Image Understanding, Elsevier, 92, 236–264, (2003)
6. Mountney, P., Lo, B., Thiemjarus, S., Stoyanov, D. and Yang, G.Z.: A Probabilistic Framework for Tracking Deformable Soft Tissue in Minimally Invasive Surgery, MICCAI, Springer, 4792, 34, (2007)
7. Zass, R. and Shashua, A.: Probabilistic graph and hypergraph matching, CVPR, (2008)
8. Torresani, L., Kolmogorov, V., and Rother, C.: Feature Correspondence via Graph Matching: Models and Global Optimization, ECCV, 596–609, (2008)
9. Leordeanu, M. and Hebert, M.: A spectral technique for correspondence problems using pairwise constraints, ICCV, 1482–1489, (2005)
10. Caetano, T.; Cheng, L.; Le, Q. and Smola, A.: Learning graph matching, ICCV, (2007)
11. Li, S.Z.: A Markov random field model for object matching under contextual constraints, CVPR, 866–869, (1994)
12. Giannarou, S., Visentini-Scarzanella, M., Yang, G.Z.: Affine-Invariant Anisotropic Detector For Soft Tissue Tracking in Minimally Invasive Surgery, ISBI, (2009)
13. Pearl, J.: Probabilistic Reasoning, San Francisco, CA: Morgan Kaufmann, (1988)
14. Mikolajczyk, K. and Schmid, C.: A Performance Evaluation of Local Descriptors, PAMI, 1615-1630, (2005)