

An EM / E-MRF Strategy for Underwater Navigation

^{1,2}Rustam Stolkin, ¹Mark Hodgetts, ²Alastair Greig
Sira Limited, UK

²Department of Mechanical Engineering

University College London, UK

Rustam.Stolkin@ptp.sira.co.uk

Mark.Hodgetts@sira.co.uk

a_greig@meng.ucl.ac.uk

Abstract

This paper addresses the problem of interpreting underwater image sequences under conditions of extremely poor visibility. The human visual system can often make correct interpretations of images that are of such poor quality that they contain insufficient explicit information to do so. We assert that this can only be achieved by utilising prior knowledge of the scene in several forms.

By modifying an Extended Markov Random Field technique, we show how image data can be dynamically combined with *expectations* of that data during image segmentation. Furthermore, we show how the interpretations of scene content and camera position can be mutually improved within an Expectation-Maximisation framework.

1 Introduction

An important problem in vision based robotic navigation, is that of recovering the position and orientation of a vehicle mounted camera relative to some recognisable object in the field of view. At any given instant, it is often possible to estimate the camera position from the recent history of the vehicle's motion. If the robot possesses a model of an important object in the scene, it is possible to refine this position estimate by firstly distinguishing that object from the background, and secondly by matching the model to the segmented image.

The second of these two problems has generated considerable interest, and a variety of strategies have been proposed for its solution. Often, however, these solutions depend on the acquisition of high quality images under constrained conditions. Typically these include constant, uniform lighting and excellent visibility.

Comparatively little attention has been given to the more difficult problem of recognising an object in the adverse conditions of the real world. In contrast we know that the human visual system is extremely robust, even in conditions of dynamic, non-uniform lighting, poor visibility and occlusion. Further, there are many instances where the human visual system can correctly interpret an image, even when that image possesses *insufficient information to enable such interpretation*. It is our belief that such a system can only function by making use of some prior knowledge of the scene.

Our algorithm (figure 1) estimates the current camera position from the recent vehicle motion using a predictive filter. A predicted (and segmented) image is then generated by projecting a 3-D model of the environment onto an image plane at the estimated position. The predicted image is used to help interpret a relatively poor visibility observed image by means of an Extended-Markov Random Field (*E-MRF*) segmentation technique. The resulting segmented image is compared with the environment model to provide a new estimate of the camera position. This improved position estimate can then be fed back into the start of the algorithm resulting in an iterative scheme which we show to be a variant of the Expectation-Maximisation (*EM*) algorithm.

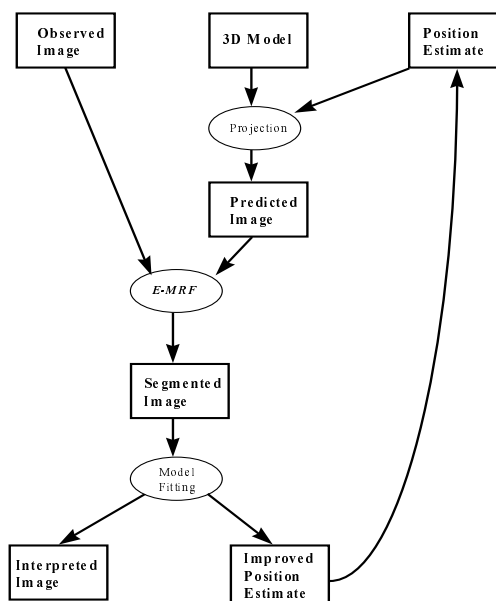


Figure 1: The *EM* / *E-MRF* strategy

2 Background

Christmas [1], [2] describes an algorithm for determining camera position by matching a projection of a known 3-D CAD model to an observed object in the image. He further suggests that the camera position in one image be used to estimate the camera position (for a moving vehicle) in the next image. Christmas demonstrates an application to autonomous navigation in a sub-sea environment, however he restricts himself to good visibility (in air) images taken of an object which is silhouetted against constant lighting in the laboratory.

Fairweather *et. al.* [3], [4], [5] tackle the problem of incomplete, poor visibility images captured under extreme, underwater conditions. They introduce the Extended-Markov Random Field (*E-MRF*) as a means of using a predicted image to help segment a poor quality observed image. They also demonstrate the superior performance of the *E-MRF* compared to both the conventional *MRF* and also a spatio-temporal *MRF* variant.

We note that Fairweather achieves correct segmentation by using a highly accurate initial estimate combined with relatively large weightings on the information taken from the predicted image. In contrast, when we apply the same method using significantly erroneous initial estimates, it fails to generate a correct segmentation although it does significantly improve upon the initial estimate. This improvement suggests an iterative scheme whereby the improved estimates are fed back into the algorithm until convergence.

Dempster [6] presents the Expectation-Maximisation (*EM*) algorithm as an iterative solution to problems where the observations can be viewed as incomplete data. The *EM* algorithm has since been presented by several authors in differing styles [7], [8], [9], [10], [11], though usually in the context of Gaussian mixture models. Neal and Hinton [7] express the algorithm in terms of calculating an expected distribution (*E-step*) for unobserved variables (in our case pixel class) in terms of observations (in our case pixel intensity) and a current estimate of parameters (in our case camera position). The Maximisation or *M-step* then re-estimates the parameters to be those with maximum likelihood. It can be shown that, with each iteration, the true likelihood improves or at least remains constant until a maximum is reached.

3 EM / E-MRF Analysis

3.1 Notation

For each image, we can define:

$\underline{\theta}$ = the (six degree-of-freedom) vector representing the true position and orientation of the camera relative to the object in the image.

$\hat{\theta}_n$ = the estimate of $\underline{\theta}$ at the beginning of the n th iteration of the algorithm.

For each pixel at a general position (i, j) in the image, we define:

$I^{i,j}$ = intensity of the pixel; $C^{i,j}$ = the true class of the pixel, either object or background.

$\hat{C}_{i,j}$ = estimate of $C^{i,j}$ from a projection of the object/environment model assuming $\hat{\theta}_n$.

It is necessary to make a distinction between $\hat{C}_{i,j}$ and $\hat{C}_{i,j}^E$ which is an adjusted estimate of $C^{i,j}$ following the segmentation stage of the algorithm. For the total image we can also define \underline{I} to represent the complete set of $I^{i,j}$ for all values of (i, j) and likewise \underline{C} , $\hat{\underline{C}}_n$ and $\hat{\underline{C}}_n^E$ respectively where n denotes the n th iteration of the algorithm. Notice that there is an obvious, geometrical one-to-one correspondence between $\underline{\theta}$ and \underline{C} , and also between $\hat{\theta}_n$ and $\hat{\underline{C}}_n$, however no such geometrical correspondence exists between $\hat{\theta}_n$ and $\hat{\underline{C}}_n^E$.

3.2 Expectation-Maximisation algorithm

Given an image containing a known object, we wish to extract the position and orientation of the camera with respect to that object. This problem may be broken down into two separate tasks. Firstly, estimating which portions of the image constitute object as opposed to background (segmentation). Secondly, estimating a camera position that best fits this segmentation. These two separate tasks at once suggest an iterative scheme that recycles the results of the second stage to re-compute the first stage and continues to perform this cycle until convergence.

We require an approach that ultimately leads us to the value of $\hat{\theta}$ that is most likely given our observations \underline{I} . Thus we seek $\hat{\theta}$ that maximises the probability $P(\underline{I} | \hat{\theta})$ or (Dempster [6], Neal and Hinton [7]) maximises a log likelihood function:

$$L(\hat{\theta}) = \log P(\underline{I} | \hat{\theta}) \quad (1)$$

Such a problem lends itself to solution by a variant of the Expectation-Maximisation (EM) algorithm which we define as follows:

E-step: Given an initial estimate $\hat{\theta}_n$ (at the start of the n th EM iteration) and corresponding set of class labels $\hat{\underline{C}}_n$, compute the expected probability:

$$E\{ P(\hat{\underline{C}}_n^E | \underline{I}; \hat{\theta}_n) \} \quad (2)$$

M-step: Compute a new position estimate $\hat{\theta}_{n+1}$ which maximises this expectation.

Within the *M-step*, our approach is firstly to segment the image, producing a new estimate of the pixel classes $\hat{\underline{C}}_n^E$ with maximum likelihood. Secondly we compute a new position estimate which best fits the pixel classes so chosen. In the following two sections we examine these stages in greater detail.

3.3 E-step

In order to compute (2) we assume:

$$E\{ P(\hat{C}_n^E | I; \underline{\theta}) \} = P(\hat{C}_n^E | I; \hat{\underline{\theta}}_n) \quad (3)$$

since $\hat{\underline{\theta}}_n$ is currently the expected value of $\underline{\theta}$. It then follows (from Bayes's law) that:

$$P(\hat{C}_n^E | I; \hat{\underline{\theta}}_n) \propto P(I | \hat{C}_n^E; \hat{\underline{\theta}}_n) \times P(\hat{C}_n^E; \hat{\underline{\theta}}_n) \quad (4)$$

3.3.1 The Extended-Markov Random Field

To compute the right hand side of equation (4) we modify a segmentation technique, first developed by Fairweather *et. al.* [3], [4], [5]. The expression consists of two parts. Fairweather calculates the first of these by initially segmenting the image off-line and using the result to generate histograms approximating the class conditional distributions $p(I \cdot | C_{i,j})$. In contrast, we have automated this process by making the approximation:

$$p(I \cdot | \hat{C}_{i,j}^E; \hat{\underline{\theta}}) \approx p(I \cdot | \hat{C}_{i,j}; \hat{\underline{\theta}}) \quad (5)$$

(The validity of this assumption depends on how closely $\hat{\underline{\theta}}$ approximates $\underline{\theta}$).

It is now comparatively easy to generate histograms of I for each value of \hat{C}_n to which we then fit normal distributions.

The second part of the expression is found with an extended Markov Random Field approach. Conventional *MRF* approaches first perform an initial segmentation by simply thresholding the image, using discriminator values derived from the normal distributions of equation (5), to give class estimates \hat{C}_n^T (i.e. the thresholding assigns a class $\hat{C}_{i,j}^T$ to every pixel (i, j)). The probability of correctness of the new classification $\hat{C}_{i,j}^E$ of a pixel (i, j) is then related to the number of that pixel's nearest neighbours that are classified similarly:

$$P(\hat{C}_{i,j}^E) \equiv P(\hat{C}_{i,j}^E; \hat{C}_{i+m, j+n}^T) \quad (6)$$

Where $(m, n \in k)$ defines a neighbourhood or clique around the pixel (i, j) . Fairweather *et. al.* extend the idea of Markov dependency by including the corresponding estimate $\hat{C}_{i,j}$ as part of this clique, thus incorporating extra (prior) knowledge of the scene into this stage of the algorithm:

$$P(\hat{C}_{i,j}^E) \equiv P(\hat{C}_{i,j}^E; \hat{C}_{i+m, j+n}^T, \hat{C}_{i,j}) \quad (7)$$

The right hand expression of equation (7) is usually expressed in the form of an exponential known as a Gibbs random field:

$$P(\hat{C}_{i,j}^E) = \frac{e^{-U_{i,j}}}{Z} \quad (8)$$

where Z is included as a normalising constant to prevent equation (8) returning probabilities greater than one.

The exponential part of this equation consists of weighted components:

$$U_{i,j} = \sum_{m, n \in k} S_1 [J(\hat{C}_{i,j}^E, \hat{C}_{i+m, j+n}^T)] + S_2 [J(\hat{C}_{i,j}^E, \hat{C}_{i,j})] \quad (9)$$

Here, S_1 and S_2 are weighting constants and J is a special function defined as:

$$J(a, b) = \begin{cases} -1 & \text{if } a = b \\ 0 & \text{if } a \neq b \end{cases} \quad (10)$$

3.4 M-step

3.4.1 Optimisation by Iterated Conditional Modes

The purpose of the *M-step* is to choose a new estimate of camera position and orientation $\hat{\theta}_{n+1}$ that maximises expression (2). We take the approach of firstly finding the optimal segmentation \hat{C}_n^E that maximises expression (2) given $\hat{\theta}_n$. We then re-estimate the camera position as $\hat{\theta}_{n+1}$ which best fits the chosen pixel classes of \hat{C}_n^E .

Optimal segmentation of the image is achieved by choosing a set of classes that maximises the right hand side of equation (4). The first part of this expression is approximated with normal distributions in equation (5) which have an exponential form. The second part of the expression is expressed as a Gibb's random field, equation (8), which also is of exponential form. Hence it is possible to combine the two into a single expression by rewriting the exponent (U part) of equation (8) to include the exponential part of the normal distributions:

i.e. U is rewritten as:

$$\sum_{i,j} S_1 [J(\hat{C}_{i,j}^E, \hat{C}_{i,j}^E)] + S_2 [J(\hat{C}_{i,j}^E, \hat{C}_{i,j}^E)] - \frac{1}{2} \log(\sigma_{c_{i,j}}^2) + \frac{(I_{i,j} - \mu_{c_{i,j}})^2}{2\sigma_{c_{i,j}}^2} \quad (11)$$

where $\sigma_{c_{i,j}}^2$ and $\mu_{c_{i,j}}$ are the variance and mean of the class conditional distribution of intensities that corresponds to the choice of $\hat{C}_{i,j}^E$ currently being considered for pixel (i, j) .

It is thus possible to achieve optimal segmentation by finding the set of classes \hat{C}_n^E that minimises the sum of all values of U over every pixel in the image.

$$\text{i.e. minimise } \sum_{\text{all } i,j} U_{i,j} \quad (12)$$

Since this actually corresponds to maximising the logarithm of the right hand side of equation (4), we consider our approach to be consistent with that of Dempster [6] and also Neal and Hinton [7] in the sense of choosing a new value of $\hat{\theta}_n$ that maximises a log likelihood function, equation (1).

There is no obvious way of choosing optimum values for the constants S_1 and S_2 . In many random field applications, such constants are chosen by trial and error. In this case we have chosen values such that the contribution from the predicted image term is approximately half that of the spatial random field and class conditional terms.

Expression (12) contains many variables since it is necessary to consider all possible classes over all pixels in the image. It is therefore not possible to search the space exhaustively to locate this minimum. We know of two established methods [12], [13] for minimising such an expression. *Simulated Annealing* is a stochastic relaxation technique which simulates the way that complex systems in the natural world achieve global energy minima. It is a powerful but relatively slow technique. In preference, we choose the method of *Iterated Conditional Modes*. This method enables decomposition of the expression (12) so that it applies to single pixels only. This method runs considerably faster, though at a risk of converging on local minima. Neal and Hinton [7] note that this kind of estimate, which improves likelihood without necessarily maximising it, will still always result in the true likelihood improving as well. Dempster [6] refers to such variants as "generalised *EM (GEM)*" algorithms.

3.4.2 Model best fitting by gradient ascent

The final stage of the *M-step* involves finding a new estimate of camera position and orientation $\hat{\theta}_{n+1}$ that best fits the new pixel classes that were chosen in the segmentation stage. We define a likelihood function for $\hat{\theta}_{n+1}$ that is based on correlation between \hat{C}_{n+1} (a function of $\hat{\theta}_{n+1}$) and \hat{C}_n . In order to maximise this likelihood function, we suggest an iterative method of gradient ascent.

4 Results

As an example, we apply the algorithm of section 3 to the problem of underwater navigation of an ROV (Remote Operated Vehicle). These unmanned submarine vehicles are used in the inspection of offshore oil-rig structures. The images used feature a scale model of a typical off-shore structure, fabricated from welded steel tubing. The images were captured by a “roving eyeball” type mini ROV at night. The only illumination is from lights mounted on the ROV itself.

Genuine underwater images because we wish to test the algorithm under the harshest visual conditions possible. The sub-sea environment is particularly challenging to machine vision systems because it contains a variety of mechanisms by which images are degraded [14], [15], [16]. Water molecules absorb light at most wavelengths and suspended particles cause scattering. Both of these mechanisms lead to severe attenuation with range. Particles larger than the wavelength of visible light cause partial occlusion. The complete absence of background lighting requires the use of vehicle mounted spotlights. These project a circular area of brightness which often saturates

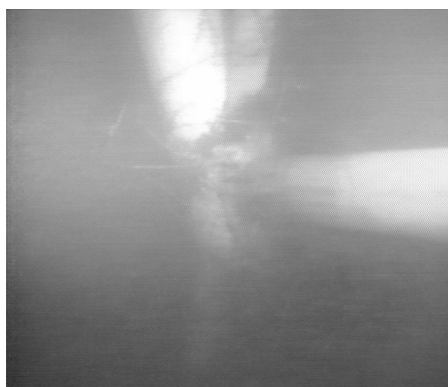


Figure 2: A typical underwater image

the camera. In contrast, areas outside of these bright patches are plunged into darkness. This uneven lighting is also constantly changing with the motion of the vehicle relative to an environment of varying reflectivity. The problem of machine perception in such a hostile environment has so far attracted comparatively little attention from researchers.

4.1 First iteration of the EM algorithm

4.1.1 Initial position estimate

The initial estimate of camera position is generated from the results of the three previous images using a simple second order predictive filter to model the motion of the ROV. This estimated position is then used to project a predicted (and segmented) image. Notice that this initial estimate is clearly and significantly erroneous when compared to the observed image (figure 3).



Figure 3: Predicted and observed images

4.1.2 Thresholding

The current position estimate is combined with the observed image, to estimate class conditional probability distributions, equation (6). The intersection of these distributions (figure 4) defines a discriminating value which is used to threshold the image (figure 5).

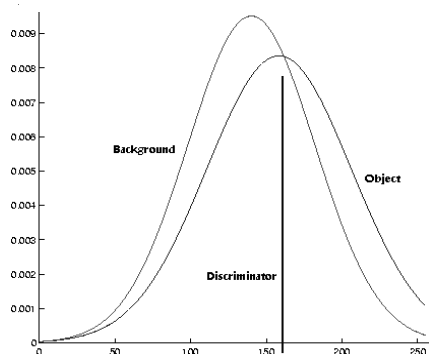


Figure 4: Class distributions

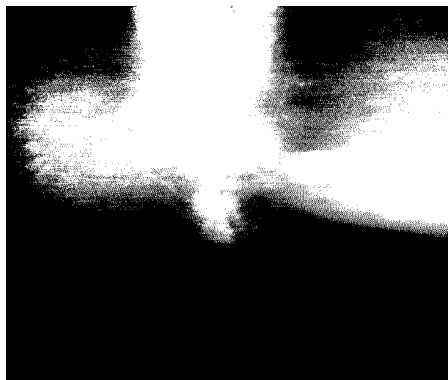


Figure 5: Thresholded image

4.1.3 E-MRF

The classes of both the neighbourhood clique and the corresponding pixel in the predicted image are used to choose the most likely class for each pixel using the *E-MRF* segmentation technique (figure 6).

The *E-MRF* extends Markov dependency to include the class of the corresponding pixel in the predicted image. To illustrate the superior performance of the *E-MRF*, we compare it here to the segmentation achieved with a conventional *MRF* (figure 7) in which Markov dependency is limited to neighbouring pixels. The conventional *MRF* incorrectly highlights an area to the left of the object. The orientation of the vertical part is incorrect and the right hand part is over-enlarged and heavily distorted.

Notice that both schemes have failed to locate the portion of object in the lower half of the picture. This is not surprising since this part of the object is entirely invisible even to the human eye. In fact the observed image contains *no information* about this part of the object. The algorithm is nevertheless still able to locate this lower portion of object by fitting its model to the segmented image (figure 8). It is interesting to note that humans who have prior knowledge of the object structure are also able to locate the lower portion in a similar fashion.



Figure 6: E-MRF segmentation



Figure 7: Conventional MRF segmentation

4.1.4 Extraction of camera position

Successive projections of the object model are best fitted (figure 8) to the segmented image by means of gradient ascent (figure 9). This camera position is used to project a new predicted image for the next iteration of the EM algorithm. Notice that this new predicted image does describe the lower part of the object even though this information does not exist in the observed image.

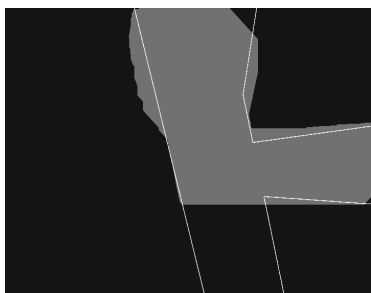


Figure 8: Model best fit

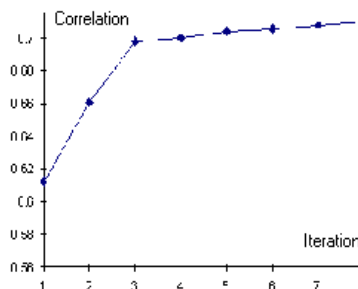


Figure 9: Improvement of correlation with gradient ascent

4.2 Successive EM iterations

In this example, the EM algorithm is seen to converge after seven iterations (figure 10). With each iteration, the algorithm produces a visible improvement on the previous position estimate (figure 11). The rate of improvement diminishes with successive iterations.



Figure 10: Final estimate

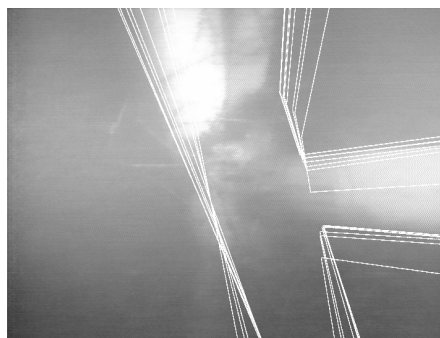


Figure 11: Improvement with successive iterations

During successive iterations, the algorithm re-estimates the class conditional probability distributions (object and background). The distribution for “background” sharpens and pulls to the left (figure 12) whereas the distribution for “object” flattens and pulls to the right (figure 13). The algorithm progressively learns to distinguish between object and background by learning that “background” is consistently dark whereas “object” is brighter but with greater intensity variation.

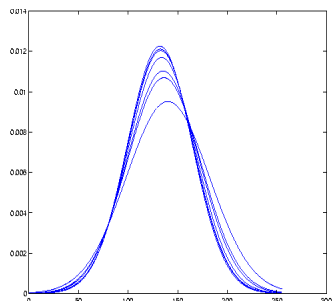


Figure 12: Improved estimate of background distribution with successive iterations

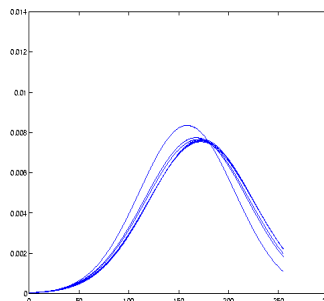


Figure 13: Improved estimate of object distribution with successive iterations

5 Performance

5.1 Accuracy

It is difficult to quantify the accuracy of the algorithm since the “ground truth” of real underwater images is difficult to measure and is rarely known. Experimental work in which artificial poor visibility scenes were modeled using computer graphics techniques [17], suggests that with extremely poor images, a very large initial error of approximately 30% on all ranges and orientation angles can be reduced to approximately 10% or less in most directions. Clearer images and smaller initial errors often result in perfect segmentation.

5.2 Comparison with Conventional Techniques

It is clear that the extended MRF with Expectation Maximisation outperforms conventional MRF segmentation as shown in figures 6 and 7 above. What is less clear is how this algorithm compares with other techniques when applied to a very low quality image. To this end we have attempted to interpret the scene using a conventional model-driven approach by fitting a CAD model of the underwater structure directly to the grey scale image. We define a goodness of fit measure that is designed to maximise the class separation (difference of means) whilst minimising the class variances:

$$\text{“Goodness of fit”} = \frac{(\mu_{object} - \mu_{background})^2}{\sigma_{object}^2 + \sigma_{background}^2} \quad (13)$$

This measure was maximised by gradient ascent. The results (figure 14) are visibly poor compared with those of the *EM / E-MRF* algorithm, given the same initial estimate (figure 3):



Figure 14: Result of fitting the model directly to the grey-scale image.



Figure 15: Superior performance of the *EM / E-MRF* algorithm.

6 Conclusions

We have described a new algorithm for the interpretation of poor visibility images. The algorithm is shown to segment a very poor visibility underwater image and significantly improve an erroneous initial estimate of camera position and orientation relative to an off-shore structure in the image. The results demonstrate a clear improvement over both conventional *MRF* and model based scene interpretation.

The algorithm allows machine vision systems to make use of prior knowledge of their environment in several novel ways. Firstly, the predicted image is used to automatically update the probability density functions required for *MRF* segmentation. Secondly, the predicted class of each pixel is introduced within an extended *MRF* model to enable image segmentation to be both data and expectation driven. Thirdly, our estimates of image interpretation and camera position are mutually refined within an Expectation-Maximisation framework.

The use of prior knowledge has enabled the algorithm to interpret parts of the image which contain little or no useful information, producing similar interpretations to those arrived at intuitively by human observers. The algorithm learns about its environment with each successive iteration and adjusts the relative contributions of the predicted and observed information by responding to the visibility conditions both at any given moment and in any given portion of the image. Future work will investigate the use of more sophisticated tracking algorithms [18] and the incorporation of underwater lighting models.

Acknowledgements

This research was undertaken within the Postgraduate Training Partnership established between Sira Ltd. and University College London. Postgraduate Training Partnerships are a joint initiative of the Department of Trade and Industry and the Engineering and Physical Sciences Research Council. They are aimed at fostering closer links between the science base, industrial research and industry.

References

- [1] W.J.Christmas, J.Kittler, M.Petrou. Error propagation for 2D-to-3D matching with application to underwater navigation. In *Proc 7th British Machine Vision Conference*, pages 555-564, 1996.
- [2] W.J.Christmas. Structural matching in computer vision using probabilistic reasoning. *PhD Thesis*, 1995.
- [3] A.J.R.Fairweather, M.A.Hodgetts, A.R.Greig. Robust Interpretation of Underwater Image Sequences. *Image Processing and its Applications* pages 660-664, 1997.
- [4] M.A.Hodgetts, A.R.Greig, A.Fairweather. Underwater Imaging Using Markov Random Fields with Feed Forward Prediction. *Journal of the Society for Underwater Technology*, Vol. 23, No.4 pages 157-167, 1999.
- [5] A.J.R.Fairweather. Robust Interpretation of Underwater Image Sequences. *PhD Thesis, UCL* 1997.
- [6] A.P.Dempster, N.M.Laird, D.B.Rubin. Maximum Likelihood from Incomplete Data via the EM Algorithm. *J. R. Stat. Soc.*, B(39) pages 1-38, 1977.
- [7] R.M.Neal,G.E.Hinton. A New View of the EM Algorithm that Justifies Incremental and Other Variants. *Biometrika*, 1993.
- [8] T.F.Cootes, C.J.Taylor. A Mixture Model for Representing Shape Variation. *Proc 8th British Machine Vision Conference*, pages 110-119, 1997.
- [9] B.North, A.Blake. Using Expectation-Maximisation to Learn Dynamical Models from Visual Data. *Proc 8th British Machine Vision Conference*, pages 669-679, 1997.
- [10] C.M.Bishop. *Neural Networks for Pattern Recognition*. Oxford University Press, 1995.
- [11] B.D.Ripley. *Pattern Recognition and Neural Networks*. Cambridge University Press, 1996.
- [12] R.C.Dubes, A.K.Jain, S.G.Nadabar, C.C.Chen. MRF Model-Based Algorithms For Image Segmentation. *Proceedings IEEE 10th International Conference on Pattern Recognition*, pages 808-814, 1990.
- [13] N.R.Pal, S.K.Pal. A Review on Image Segmentation Techniques. *Pattern Recognition*, Vol26, No.9, pages 1277-1294, 1993.
- [14] A.Morel. Optical Properties of Pure Water and Sea Water. *Optical Aspects of Oceanography*, pages 1-24, 1974.
- [15] Kullenburg. Observed and Computed Scattering Functions. *Optical Aspects of Oceanography*, pages 25-49, 1974.
- [16] Jerlov. *Marine Optics. Elsevier Oceanography Series 14*. Amsterdam Elsevier Scientific Publishing Company, 1976.
- [17] P.Rokita. Simulating Poor Visibility Conditions Using Image Processing. *Real-Time Imaging 3*, pages 275-281, 1997.
- [18] M.Isard, A.Blake. Contour tracking by stochastic propagation of conditional density. *Proc. European Conf. Computer Vision*, pages 343-356, 1996.