

Pel-Recursive Motion Estimation Using the Expectation-Maximization Technique and Spatial Adaptation

Vania Estrela⁽¹⁾, Luis A. Rivera⁽¹⁾, Marcos H. S. Bassani⁽²⁾

⁽¹⁾Universidade Estadual do Norte Fluminense, LCMAT-CCT, Grupo de Computação Científica
Av. Alberto Lamego 2000, Campos dos Goytacazes,
RJ, Brasil, CEP 28013-600
{vaniave,rivera}@uenf.br

⁽²⁾Departamento de Engenharia Mecânica, DEPES, CEFET-RJ
Av. Maracanã 229, Rio de Janeiro, RJ, BRAZIL, CEP 20271-110
marcos@cefet-rj.br

ABSTRACT

Pel-recursive motion estimation is a well-established approach. However, in the presence of noise, it becomes an ill-posed problem that requires regularization. In this paper, motion vectors are estimated in an iterative fashion by means of the Expectation-Maximization (EM) algorithm and a Gaussian data model. Our proposed algorithm also utilizes the local image properties of the scene to improve the motion vector estimates following a spatially adaptive approach. Numerical experiments are presented that demonstrate the merits of our methods.

Keywords

Motion estimation, EM algorithm, pel-recursive, maximum likelihood.

1. INTRODUCTION

Motion estimation is very important in multimedia video processing applications. For example, in video coding, the estimated motion is used to reduce the transmission bandwidth. The evolution of an image sequence motion field can also help other image processing tasks in multimedia applications such as analysis, recognition, tracking, restoration, collision avoidance and segmentation of objects [Jain89].

In coding applications, a block-based approach [Tek95] is often used for interpolation of lost information between key frames. The fixed rectangular partitioning of the image used by some block-based approaches often separates visually meaningful image features. If the components of an important feature are assigned different motion vectors, then the interpolated image will suffer from

annoying artifacts.

Pel-recursive schemes [Brailean95, Estrela03, Jain89] can theoretically overcome some of the limitations associated with blocks by assigning a unique motion vector to each pixel. Intermediate frames are then constructed by resampling the image at locations determined by linear interpolation of the motion vectors. The pel-recursive approach can also manage motion with sub-pixel accuracy. However, its original formulation was deterministic. The update of the motion estimate was based on the minimization of the displaced frame difference (DFD) at a pixel. In the absence of additional assumptions about the pixel motion, this estimation problem becomes ill-posed because of the following problems: a) occlusion; b) the solution to the 2D motion estimation problem is not unique; and c) the solution does not continuously depend on the data due to the fact that motion estimation is highly sensitive to the presence of observation noise in video images.

In this work, we plan to use the MAP estimate to find the update of the motion vector from our observation model by means of the Expectation-Maximization technique. The main advantages of the EM method is that the final algorithm deals with

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

WSCG'2004, February 2-6, 2003, Plzen, Czech Republic.
Copyright UNION Agency – Science Press

closed-form expressions and it does not require the use of optimization techniques.

We organized this work as follows. Section 2 provides some necessary background on the pel-recursive motion estimation problem. Section 3 introduces our spatially adaptive approach. Section 4 describes the EM framework. Section 5 defines the metrics used to evaluate our results. Section 6 describes some implementation aspects and the experiments used to access the performance of our proposed algorithm. Finally, Section 7 has some conclusions.

2. PEL-RECURSIVE DISPLACEMENT ESTIMATION

The displacement of each picture element (pel or pixel) in each frame forms the displacement vector field (DVF) and its estimation can be done using at least two successive frames. The DVF is the 2D motion resulting from the apparent motion of the image brightness. A vector is assigned to each point in the image.

A pixel belongs to a moving area if its intensity has changed between consecutive frames. Hence, our goal is to find the corresponding intensity value $I_k(\mathbf{r})$ of the k -th frame at location $\mathbf{r} = [x, y]^T$, and the corresponding displacement vector (DV) at the working point \mathbf{r} in the current frame, that is $\mathbf{d}(\mathbf{r}) = [d_x, d_y]^T$. Pel-recursive algorithms minimize the DFD function in a small area containing the working point assuming constant image intensity along the motion trajectory. The DFD is defined by

$$\Delta(\mathbf{r}; \mathbf{d}(\mathbf{r})) = I_k(\mathbf{r}) - I_{k-1}(\mathbf{r} - \mathbf{d}(\mathbf{r})) \quad (1)$$

and the perfect registration of frames will result in $I_k(\mathbf{r}) = I_{k-1}(\mathbf{r} - \mathbf{d}(\mathbf{r}))$. The DFD represents the error due to the nonlinear temporal prediction of the intensity field through the DV. The relationship between the DVF and the intensity field is nonlinear. An estimate of $\mathbf{d}(\mathbf{r})$, is obtained by directly minimizing $\Delta(\mathbf{r}, \mathbf{d}(\mathbf{r}))$ or by determining a linear relationship between these two variables through some model. This is accomplished by means of a Taylor series expansion of $I_{k-1}(\mathbf{r} - \mathbf{d}(\mathbf{r}))$ about the location $(\mathbf{r} - \mathbf{d}^i(\mathbf{r}))$, where $\mathbf{d}^i(\mathbf{r})$ represents a prediction of $\mathbf{d}(\mathbf{r})$ in the i -th step. This results in

$$\Delta(\mathbf{r}, \mathbf{r} - \mathbf{d}^i(\mathbf{r})) = -\mathbf{u}^T \nabla I_{k-1}(\mathbf{r} - \mathbf{d}^i(\mathbf{r})) + e(\mathbf{r}, \mathbf{d}(\mathbf{r})), \quad (2)$$

where the displacement update vector $\mathbf{u} = [u_x, u_y]^T = \mathbf{d}(\mathbf{r}) - \mathbf{d}^i(\mathbf{r})$; $e(\mathbf{r}, \mathbf{d}(\mathbf{r}))$ represents the error resulting from the truncation of the higher order terms (linearization error); and $\nabla = [\partial/\partial x, \partial/\partial y]^T$ represents the spatial gradient operator. Applying (2) to all points in a neighborhood \mathcal{R} gives

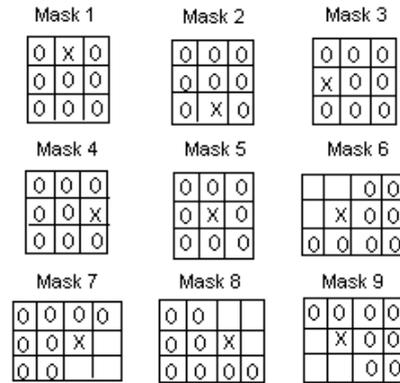
$$\mathbf{z} = \mathbf{G}\mathbf{u} + \mathbf{n}, \quad (3)$$

where the temporal gradients $\Delta(\mathbf{r}, \mathbf{r} - \mathbf{d}^i(\mathbf{r}))$, from (2), have been stacked to form the $N \times 1$ observation vector \mathbf{z} containing DFD information on all the pixels in a neighborhood \mathcal{R} , the $N \times 2$ matrix \mathbf{G} is obtained by stacking the spatial gradient operators at each observation, and the error terms have formed the $N \times 1$ noise vector \mathbf{n} which is assumed Gaussian with probability density function (pdf) $\mathbf{n} \sim \mathcal{N}(\mathbf{0}, \sigma_n^2 \mathbf{I})$. Each row of \mathbf{G} has entries $[g_{xi}, g_{yi}]^T$, with $i = 1, \dots, N$. The spatial gradients of I_{k-1} are calculated through a bilinear interpolation scheme [Brailean95].

3. SPATIAL ADAPTION

Aiming to improve the estimates given by the pel-recursive algorithm, we introduced an adaptive scheme for determining the optimal shape of the neighborhood of pixels overdetermined system of equations given by (3).

Errors can be caused by the basic underlying assumption of uniform motion inside \mathcal{R} (the smoothness constraint), by not grouping pixels adequately, and by the way gradient vectors are estimated, among other things. Since it is known that in a noiseless image not containing pixels with constant intensity, most errors, when estimating motion, occur close to motion boundaries. We propose a hypothesis testing approach to determine the best neighborhood shape for a given pixel. The masks in Figure 1 show the geometries of the neighborhoods used, where “x” indicates the position of the current pixel and “0” is a neighboring pixel. We pick up the neighborhood from the finite set of templates shown in Figure 1, according to the smallest $|DFD|$ criterion, in an attempt to adapt the model to local features associated to motion boundaries. The chosen mask is the one that better satisfies the smoothness constraint, capturing the motion behaviour around the current pixel, and, therefore accounting for motion discontinuities.



X Current pixel 0 Neighboring pixel

Figure 1. Neighborhood geometries.

4. THE EM APPROACH

The EM algorithm is a general numerical technique which can be used to determine the maximum likelihood estimate (MLE) of a set of parameters. It can be employed for identification of model and distribution parameters simultaneously. The parameters can be estimated iteratively even in situations where some variables cannot be observed.

4.1 The Ordinary Maximum Likelihood Estimate and Its Limitations

The vector \mathbf{u} of (3) can be computed using the MAP estimate [Kay93] as

$$\hat{\mathbf{u}}_{MAP} = \mathbf{A}_u \mathbf{G}^T (\mathbf{G} \mathbf{A}_u \mathbf{G}^T + \mathbf{A}_n)^{-1} \mathbf{z}, \quad (4)$$

where \mathbf{A}_u and \mathbf{A}_n are, respectively, the covariance matrices of the update vector and of the noise/linearization error term.

The estimate in (4) was derived assuming that \mathbf{u} and \mathbf{n} are zero mean, and uncorrelated. The maximum a posteriori (MAP) estimate for the linear observation model in (3) is also the MAP estimate assuming a Gaussian prior on \mathbf{u} and \mathbf{n} .

In this work, we resort to the maximum likelihood (ML) estimation of the second-order statistics \mathbf{A}_u and \mathbf{A}_n of the model. The calculation of the ML estimates of \mathbf{A}_u and \mathbf{A}_n is done iteratively by means of the Expectation-Maximization (EM) algorithm formalized by Dempster et al [Dem77] which has been used in a variety of applications as referenced in [Kat91]. The EM algorithm was previously used for motion estimation in [Fan96]. However, this work estimates parameters of affine models and it uses a block-matching framework.

For the model described by (3), let us assume that the update vector \mathbf{u} and the noise \mathbf{n} are normally-distributed with mean zero, and uncorrelated, with covariance matrices equal to \mathbf{A}_u and \mathbf{A}_n respectively. Then, the pdf of \mathbf{z} is given by

$$f_z(\mathbf{z}) = \left| 2\pi(\mathbf{G} \mathbf{A}_u \mathbf{G}^H + \mathbf{A}_n) \right|^{-\frac{1}{2}} \times \exp \left\{ -\frac{1}{2} \mathbf{z}^H (\mathbf{G} \mathbf{A}_u \mathbf{G}^H + \mathbf{A}_n)^{-1} \mathbf{z} \right\}.$$

where \mathbf{G}^H and \mathbf{z}^H are the Hermitian conjugates of \mathbf{G} and \mathbf{z} [Kat91].

Let us take $\Phi = \{\mathbf{A}_u, \mathbf{A}_n\}$ as the parameter set to be estimated. The pdf of \mathbf{z} can be considered a continuous and differentiable function of Φ . This dependency can be better captured by the notation $f_z(\mathbf{z}; \Phi)$. Furthermore, we can assume that the additive noise is white, with covariance matrix $\mathbf{A}_n = \sigma_n^2 \mathbf{I}$, where \mathbf{I} is an $N \times N$ identity matrix.

The ML estimation of Φ is the Φ_{ML} that maximizes the logarithm of the likelihood function of $f_z(\mathbf{z}; \Phi)$, which is given by

$$\begin{aligned} \Phi_{ML} &= \arg \{ \max_{\Phi} f_z(\mathbf{z}; \Phi) \} \\ &= \arg \{ \max_{\Phi} \log f_z(\mathbf{z}; \Phi) \}. \end{aligned} \quad (5)$$

Combining (4) and (5), we conclude that the maximization of the log-likelihood function is equivalent to minimizing the function $L_o(\Phi)$ where

$$L_o(\Phi) = \log \left| \mathbf{G} \mathbf{A}_u \mathbf{G}^H + \mathbf{A}_n \right| + \mathbf{z}^H (\mathbf{G} \mathbf{A}_u \mathbf{G}^H + \mathbf{A}_n)^{-1} \mathbf{z}. \quad (6)$$

This function is nonlinear and non-unimodal with respect to Φ . Analytical solutions for (6) are out of question. Thus, we have to resort to optimization techniques. Nevertheless, since $L_o(\Phi)$ is not unimodal, convergence and local extrema of iterative optimization methods are serious problems. The EM algorithm arises as an alternative for estimating Φ , since its convergence to a local maximum is guaranteed.

4.2 Problem Formulation According to the EM Framework

We assume that the update vector \mathbf{u} is normally-distributed with mean zero and covariance matrix equal to \mathbf{A}_u , and that $\mathbf{n} \sim N(\mathbf{0}, \sigma_n^2 \mathbf{I})$ is the additive Gaussian noise.

If we choose $\mathbf{x} = [\mathbf{u} \ \mathbf{n}]^T$ as the complete data where $\mathbf{x} \in \mathbb{R}^{(N+2) \times 1}$, then \mathbf{x} will be normally distributed with diagonal covariance matrix \mathbf{A}_x equal to

$$\mathbf{A}_x = \begin{bmatrix} \mathbf{A}_u & \mathbf{0} \\ \mathbf{0} & \mathbf{A}_n \end{bmatrix} \quad \text{and} \quad \mathbf{A}_x^{-1} = \begin{bmatrix} \mathbf{A}_u^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{A}_n^{-1} \end{bmatrix},$$

with this, (3) can be re-written as

$$\mathbf{z} = [\mathbf{G}; \mathbf{I}] \mathbf{x} = \mathbf{H} \mathbf{x}, \quad \text{for } \mathbf{H} = [\mathbf{G}; \mathbf{I}].$$

The foundation of the EM algorithm is the maximization of the expectation of $\log \{f_x(\mathbf{x}; \Phi)\}$, given the incomplete observed data \mathbf{z} and the current estimate of the parameter set Φ . $f_x(\mathbf{x}; \Phi)$ is the pdf of the complete data, and it is given by

$$f_x(\mathbf{x}; \Phi) = \left| 2\pi \mathbf{A}_x \right|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} \mathbf{x}^H \mathbf{A}_x^{-1} \mathbf{x} \right\}, \quad (7)$$

that taking the logarithm of both sides leads to

$$\begin{aligned} \log \{f_x(\mathbf{x}; \Phi)\} &= -\frac{1}{2} \log \left| 2\pi \mathbf{A}_x \right| - \frac{1}{2} \mathbf{x}^H \mathbf{A}_x^{-1} \mathbf{x} \\ &= K - \frac{1}{2} [\mathbf{A}_x + \mathbf{x}^H \mathbf{A}] \end{aligned}$$

where $\Phi = \mathbf{A}_x$ and K is a constant independent of Φ .

4.1.1 EM algorithm

We want to compute iteratively

$$\Phi^{(p+1)} = \arg\{ \max_{\Phi} Q(\Phi; \Phi^{(p)}) \}. \quad (8)$$

For that, we need to compute

$$Q(\Phi; \Phi^{(p)}) = E[\log\{ f_x(\mathbf{x}; \Phi) \} | \mathbf{z}; \Phi^{(p)}], \quad (9)$$

where $\Phi^{(p)} = [\mathbf{A}_u^{(p)}; \mathbf{A}_n^{(p)}]$ is the estimate of the parameter set $\Phi = \{\mathbf{A}_u; \mathbf{A}_n\}$ at the p -th iteration.

We compute (9) and (8) in E-step and M-step, respectively.

E-Step:

This step can be re-written in terms of the parameter set Φ by means of the relationships developed in terms of the complete data \mathbf{x} and the pdf $f(\mathbf{x}; \Phi^{(p)})$ as

$$\begin{aligned} F(\Phi; \Phi^{(p)}) &= \log |\mathbf{A}_u| + \log |\mathbf{A}_n| + \text{tr}(\mathbf{A}_u^{-1} \mathbf{A}_u^{(p)}) \\ &+ \text{tr}(\mathbf{A}_n^{-1} \mathbf{A}_n^{(p)}) + \mu_{u|z}^{(p)H} \mathbf{A}_u^{-1} \mu_{u|z}^{(p)} \\ &+ \mu_{n|z}^{(p)H} \mathbf{A}_n^{-1} \mu_{n|z}^{(p)} \end{aligned} \quad (10)$$

Our goal is to estimate the update vector \mathbf{u} . The MAP/MMSE estimate of \mathbf{u} is $\mu_{u|z}^{(p)}$ which corresponds to its conditional expectation at iteration p , given the observation \mathbf{z} . Since our observation model is linear and Gaussian statistics are assumed, we have $\hat{\mathbf{u}}_{MMSE} = \hat{\mathbf{u}}_{LMSE} = \hat{\mathbf{u}}_{MAP}$.

This estimate is obtained as a byproduct of the estimation of Φ . The statistical assumptions about \mathbf{u} and \mathbf{n} result in $\Phi = \{\sigma_1^2, \sigma_2^2, \sigma_n^2\}$. Therefore,

$$\begin{aligned} F(\Phi; \Phi^{(p)}) &= \log(\sigma_1^2) + \log(\sigma_2^2) + \\ &m[\log(\sigma_n^2)] + \frac{a_{11}^{(p)} + c_1^{2(p)}}{\sigma_1^2} + \frac{a_{22}^{(p)} + c_2^{2(p)}}{\sigma_2^2} \\ &+ \frac{1}{\sigma_n^2} [b_{11}^{(p)} + \dots + b_{mm}^{(p)} + e_1^{2(p)} + \dots + e_m^{2(p)}] \end{aligned} \quad (11)$$

where

$$\begin{aligned} \mathbf{A}^{(p)} &= [\Lambda_u^{(p)} - \Lambda_u^{(p)} \mathbf{G}^H (\mathbf{G} \Lambda_u^{(p)} \mathbf{G}^H + \Lambda_n^{(p)})^{-1} \mathbf{G} \Lambda_u^{(p)}], \\ \mathbf{A}^{(p)} &= [\Lambda_u^{(p)} - \Lambda_u^{(p)} \mathbf{G}^H (\mathbf{G} \Lambda_u^{(p)} \mathbf{G}^H + \Lambda_n^{(p)})^{-1} \mathbf{G} \Lambda_u^{(p)}], \\ \mathbf{B}^{(p)} &= [\Lambda_n^{(p)} - \Lambda_n^{(p)} (\mathbf{G} \Lambda_u^{(p)} \mathbf{G}^H + \Lambda_n^{(p)})^{-1} \Lambda_n^{(p)}], \\ \mathbf{c}^{(p)} &= \mu_{u|z}^{(p)} = [\Lambda_u^{(p)} \mathbf{G}^H (\mathbf{G} \Lambda_u^{(p)} \mathbf{G}^H + \Lambda_n^{(p)})^{-1} \mathbf{z}] \text{ and} \\ \mathbf{e}^{(p)} &= \mu_{n|z}^{(p)} = [\Lambda_n^{(p)} (\mathbf{G} \Lambda_u^{(p)} \mathbf{G}^H + \Lambda_n^{(p)})^{-1} \mathbf{z}]. \end{aligned} \quad (12)$$

The previous matrices can be simplified and written as

$$\mathbf{A}^{(p)} = \begin{bmatrix} a_{11}^{(p)} & a_{12}^{(p)} \\ a_{21}^{(p)} & a_{22}^{(p)} \end{bmatrix}, \quad \mathbf{B}^{(p)} = \begin{bmatrix} b_{11}^{(p)} & \dots & b_{1m}^{(p)} \\ \vdots & \ddots & \vdots \\ b_{m1}^{(p)} & \dots & b_{mm}^{(p)} \end{bmatrix},$$

$$\mathbf{c}^{(p)} = [c_1^{(p)} \quad c_2^{(p)}]^T \text{ and } \mathbf{e}^{(p)} = [e_1^{(p)} \quad \dots \quad e_m^{(p)}]^T.$$

M-Step:

The minimum of $F(\Phi; \Phi^{(p)})$ occurs when $\frac{\partial F(\Phi; \Phi^{(p)})}{\partial \sigma_n^2} = 0$. So,

$$\sigma_n^{2(p+1)} = \frac{1}{m} \left\{ \text{Tr}[\mathbf{B}^{(p)}] + |\mathbf{e}^{(p)}|^2 \right\}. \quad (13)$$

Applying similar procedure for σ_1^2 and σ_2^2 gives

$$\sigma_1^{2(p+1)} = a_{11}^{(p)} + c_1^{2(p)} \text{ and } \sigma_2^{2(p+1)} = a_{22}^{(p)} + c_2^{2(p)}. \quad (14)$$

Now, we are ready to state the resulting EM algorithm using multiple masks.

4.3 The EM-Based Pel Recursive Motion Estimation Algorithm

For each pixel in the current frame k , located at \mathbf{r} , do the following:

Initialize: $\mathbf{d}^0(\mathbf{r})$, $\sigma_n^{2(0)}$, $\sigma_1^{2(0)}$ and $\sigma_2^{2(0)}$;

$m = p \leftarrow 0$ (mask and iteration counter);

thresholds T, ε and ξ .

Do {

Calculate \mathbf{G}^p , \mathbf{z}^p for current mask

Calculate $\mathbf{c}^{(p)}$ from (12)

Perform M-Step

Calculate $\mathbf{d}^{p+1}(\mathbf{r}) = \mathbf{d}^p(\mathbf{r}) + \mathbf{u}^p$

$p \leftarrow p + 1$

If ($p = MAX$) then

$m \leftarrow m + 1$ (try another neighborhood)

reset $\mathbf{d}^0(\mathbf{r})$, $\sigma_n^{2(0)}$, $\sigma_1^{2(0)}$, $\sigma_2^{2(0)}$, $p \leftarrow 0$

end if

if (all masks where used), $\mathbf{d}^{p+1}(\mathbf{r}) = \mathbf{0}$

} while ($\|\Phi^{p+1} - \Phi^p\| \leq \xi$ and

$\|\mathbf{d}^{p+1}(\mathbf{r}) - \mathbf{d}^p(\mathbf{r})\| \leq \varepsilon$ and $|DFD| < T$

and $p < MAX$)

stop.

5. METRICS

This work assesses the motion field quality through the use of the four metrics [Brailean95, Estrela03] as described below.

5.1 Mean Squared Error (MSE)

Since the MSE provides an indication of the degree of correspondence between the estimates and the true value of the motion vectors, we can apply this measure to two consecutive frames of a sequence with known motion. We can evaluate the MSE in the horizontal (MSE_x) and in the vertical (MSE_y) directions as follows:

$$MSE_x = \frac{1}{RC} \sum_{\mathbf{r} \in S} [d_x(\mathbf{r}) - \hat{d}_x(\mathbf{r})]^2, \text{ and}$$

$$MSE_y = \frac{1}{RC} \sum_{\mathbf{r} \in S} [d_y(\mathbf{r}) - \hat{d}_y(\mathbf{r})]^2,$$

where S is the entire frame, \mathbf{r} represents the pixel coordinates, R and C are, respectively, the number of rows and columns in a frame, $\mathbf{d}(\mathbf{r})=(d_x(\mathbf{r}), d_y(\mathbf{r}))$ is the true deviation vector at \mathbf{r} , and $\hat{\mathbf{d}}(\mathbf{r})=(\hat{d}_x(\mathbf{r}), \hat{d}_y(\mathbf{r}))$ its estimation.

5.2 Bias

The bias gives an idea of the degree of correspondence between the estimated motion field and the original optical flow. It is defined as the average of the difference between the true DV's and their predictions, for all pixels inside a frame S , and it is defined along the x and y directions as

$$bias_x = \frac{1}{RC} \sum_{\mathbf{r} \in S} [d_x(\mathbf{r}) - \hat{d}_x(\mathbf{r})], \text{ and}$$

$$bias_y = \frac{1}{RC} \sum_{\mathbf{r} \in S} [d_y(\mathbf{r}) - \hat{d}_y(\mathbf{r})]$$

5.3 Mean-Squared Displaced Frame Difference

This metric evaluates the behavior of the average of the squared displaced frame difference (\overline{DFD}^2). It represents an assessment of the evolution of the temporal gradient as the scene evolves by looking at the squared difference between the current intensity $I_k(\mathbf{r})$ and its predicted value $I_{k-1}(\mathbf{r}-\mathbf{d}(\mathbf{r}))$. Ideally, the \overline{DFD}^2 should be zero, which means that all motion was identified correctly ($I_k(\mathbf{r})=I_{k-1}(\mathbf{r}-\mathbf{d}(\mathbf{r}))$ for all \mathbf{r} 's). In practice, we want the \overline{DFD}^2 to be as low as possible. Its is defined as

$$\overline{DFD}^2 = \frac{\sum_{k=2}^K \sum_{\mathbf{r} \in S} [I_k(\mathbf{r}) - I_{k-1}(\mathbf{r}-\mathbf{d}(\mathbf{r}))]^2}{RC(K-1)},$$

where K is the length of the image sequence.

5.3 Improvement in Motion Compensation

The improvement in motion compensation ($\overline{IMC}(dB)$) between two consecutive frames is given by

$$\overline{IMC}_k(dB) = 10 \log_{10} \left\{ \frac{\sum_{\mathbf{r} \in S} [I_k(\mathbf{r}) - I_{k-1}(\mathbf{r})]^2}{\sum_{\mathbf{r} \in S} [I_k(\mathbf{r}) - I_{k-1}(\mathbf{r}-\mathbf{d}(\mathbf{r}))]^2} \right\},$$

where S is the frame being currently analyzed. It shows the ratio in decibel (dB) between the mean-squared frame difference (\overline{FD}^2) defined by

$$\overline{FD}^2 = \frac{\sum_{\mathbf{r} \in S} [I_k(\mathbf{r}) - I_{k-1}(\mathbf{r})]^2}{RC}$$

and the \overline{DFD}^2 between frames k and $(k-1)$.

As far as the use of this metric goes, we chose to apply it to a sequence of K frames, resulting in the following equation for the average improvement in motion compensation:

$$\overline{IMC}(dB) = 10 \log_{10} \left\{ \frac{\sum_{k=2}^K \sum_{\mathbf{r} \in S} [I_k(\mathbf{r}) - I_{k-1}(\mathbf{r})]^2}{\sum_{k=2}^K \sum_{\mathbf{r} \in S} [I_k(\mathbf{r}) - I_{k-1}(\mathbf{r}-\mathbf{d}(\mathbf{r}))]^2} \right\}$$

When it comes to motion estimation, we seek algorithms that have high values of $\overline{IMC}(dB)$. If we could detect motion without any error, then the denominator of the previous expression would be zero (perfect registration of motion) and we would have $\overline{IMC}(dB) = \infty$.

6. IMPLEMENTATION

In this section, we present several experimental results illustrating the effectiveness of the EM algorithm and compare it with the Wiener filter described by

$$\hat{\mathbf{u}}_{Wiener} = \hat{\mathbf{u}}_{LMMSE} = (\mathbf{G}^T \mathbf{G} + \mu \mathbf{I})^{-1} \mathbf{G}^T \mathbf{z},$$

with $\mu=50$ in which the statistics of \mathbf{A}_i and \mathbf{A}_n remain constant for the entire frame [Bie87]. As before, the algorithms were tested on a synthetic sequence and two real video sequences: the "Foreman" and the "Mother and Daughter". The sequences are 144 x 176, 8-bit (QCIF).

6.1 Experiment 1

In this sequence, there is a moving rectangle immersed in a moving background. In order to create textures for the rectangle and its background (otherwise motion detection would not be possible), the following auto-regressive model was used:

$$I(m,n) = \frac{1}{3}[I(m,n-1)+I(m-1,n)+I(m-1,n-1)] + n_i(m,n),$$

where $i=1,2$. For the background ($i=1$), n_i is a Gaussian random variable with mean $\mu_1 = 50$ and variance $\sigma_1^2 = 49$. The rectangle ($i=2$) was generated with $\mu_2 = 100$ and variance $\sigma_2^2 = 25$. All pixels from the background move to the right, and the displacement from frame 1 to frame 2 is $\mathbf{d}_b(\mathbf{r})=(d_{bx}(\mathbf{r}), d_{by}(\mathbf{r}))=(2,0)$. The rectangle moves diagonally from frame 1 to 2 with $\mathbf{d}_r(\mathbf{r})=(d_{rx}(\mathbf{r}), d_{ry}(\mathbf{r}))=(1,2)$.

Table 1 shows the values for the MSE, bias, $\overline{IMC}(dB)$ and \overline{DFD}^2 for the estimated image brightness for frames 1 and 2 of the noiseless

"Synthetic" sequence using the following algorithms: Wiener filter, EM with one mask (EM), and EM multi-mask (EMm). All implementations using the EM technique performed better than the Wiener filter. Table 2 shows the values of the same metrics when $SNR = 20dB$.

The motion-compensated frames corresponding to the estimated DVF's for the noisy case ($SNR=20dB$) are presented in Figure 2. Visually speaking, there is an improvement in the outcome around the borders of the rectangle for the multi-mask implementation of the EMm.

6.2 Experiment 2.

Figure 3 presents the values of the \overline{IMC} for frames 31 to 40 of the "Mother and Daughter" sequence for the noiseless and noisy ($SNR=20dB$) cases, respectively, for the multi-mask algorithm EMA=EMm. If we look at the noiseless case, we will see that the algorithm EMA=EMm has very good performance. For the noisy case EMA=EMm algorithm can perform extremely well most of the time, except for two pairs of frames. Their qualitative (visual) performance can be observed in Figure 4. The superior performance of the EMA=EMm routine can be better noticed by looking at the errors resulting from motion compensation.

6.3 Experiment 3

Figure 5 presents the values of the \overline{IMC} for frames 11 to 20 of the "Foreman" sequence for the noiseless and noisy ($SNR=20dB$) cases, respectively, for the multi-mask algorithm EMm. The EM implementations outperform the Wiener filter all the time. Their qualitative performance can be observed in Figure 6 which shows the error in motion compensation. Visually speaking, we perceive that the EM seems to be more robust to the motion introduced by the camera, because the background of the frame showing the errors due to motion-compensation has less artifacts than the one generated with the Wiener filter.

7. CONCLUSIONS

For all the sequences, the EM algorithm performed better than the Wiener filter for both the one mask and multi-mask cases, regardless of the presence of noise.

The EM algorithm showed some sensitivity to the choice of initial estimates. We used more than one

initial parameter set $\Phi^0 = [\sigma_n^{2(0)}, \sigma_1^{2(0)}, \sigma_2^{2(0)}]^T$ to improve the rate of convergence, but even with this extra feature, the resulting algorithms using one mask and multiple masks were faster than the corresponding Wiener filter counterparts.

For the EM method, we have a simple algorithm that is guaranteed to converge and it does not require numerical optimization. Given that there are multiple iterations at every pixel location, the speed advantage gained by means of the EM algorithm is considerable.

The authors think this framework has great potential for applications such as video coding and image segmentation.

8. REFERENCES

- [Bie87] Biemond, J.; Looijenga, L.; Boeke, D.E. and Plompen, R. H. J. M. A Pel-Recursive Wiener-Based Displacement Estimation Algorithm, Signal Processing, 13, pp. 399-412, December 1987.
- [Brailean95] Brailean, J.C. and Katsaggelos, A. K. Simultaneous Recursive Displacement Estimation and Restoration of Noisy-Blurred Image Sequences, IEEE Trans. Im. Proc., Vol. 4, No. 9, pp. 1236-1268, September 1995.
- [Estrela03] Estrela, V.; Rivera, L.A.; Beggio, P.C. and Lopes, R.T. Regularized Pel-Recursive Motion Estimation Using Generalized Cross-Validation and Spatial Adaptation, Proc. of SIBGRAPI2003 XVI Brazilian Symposium Computer Graphic and Image Processing, pp. 331-338, Oct. 2003.
- [Dem77] Dempster, A.P.; Laird, N.M. and Rubin, D.B. Maximum Likelihood From Incomplete Data via the EM Algorithm, Ann. Roy. Statist. Soc., Vol. 39, pp. 1-38, Dec. 1977.
- [Fan96] Fan, C.M.; Namazi, N.M. and Penafiel, P.B. A New Image Motion Estimation Algorithm Based on the EM Technique, IEEE Trans. on P.A.M.I., Vol. 18, No. 3, pp. 348-352, 1996.
- [Jain89] Jain, A.K. Fundamentals of Digital Image Processing, Prentice-Hall, New Jersey, 1989.
- [Kat91] Katsaggelos, A.K. Image Restoration, Springer-Verlag, Berlin, Heidelberg, 1991.
- [Kay93] Kay, S. Fundamentals of Statistical Signal Processing: Estimation Theory, Prentice Hall, 1993.
- [Tek95] Tekalp, M., Digital Video Processing, Prentice Hall, 1995.

	Wiener	EM	EMm
MSE_x	0.1548	0.1436	0.1311
MSE_y	0.0740	0.0681	0.0572
$bias_x$	0.0610	0.0571	0.0527
$bias_y$	-0.0294	-0.0189	-0.0180
$\overline{IMC} (dB)$	19.46	19.93	20.47
\overline{DFD}^2	4.16	4.08	3.26

Table 1. Comparison between EM implementations and the Wiener filter. $SNR = \infty$.

	Wiener	EM	EMm
MSE_x	0.2563	0.2401	0.2323
MSE_y	0.1273	0.1249	0.1226
$bias_x$	0.0908	0.0842	0.0822
$bias_y$	-0.0560	-0.0558	-0.0503
$\overline{IMC}(dB)$	14.74	15.09	15.45
\overline{DFD}^2	12.24	10.91	10.03

Table 2. Comparison between EM implementations and the Wiener filter. $SNR = 20dB$.

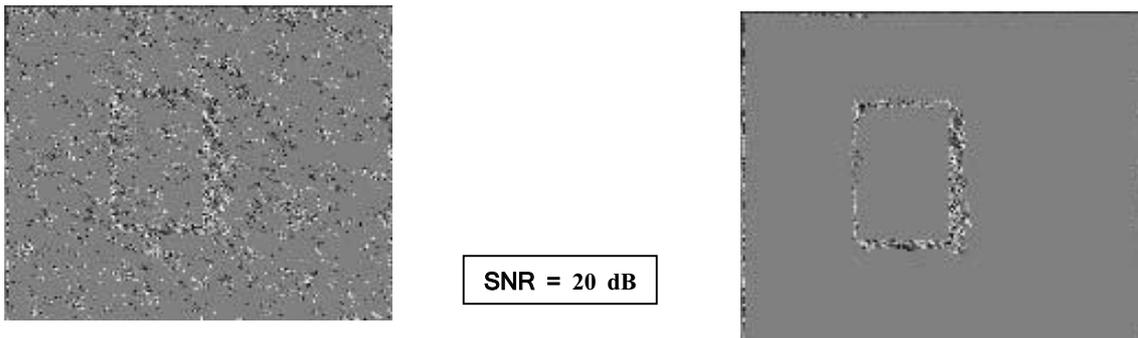


Figure 2. Motion-compensation errors for the Wiener (left) and the EMm (right) algorithms.

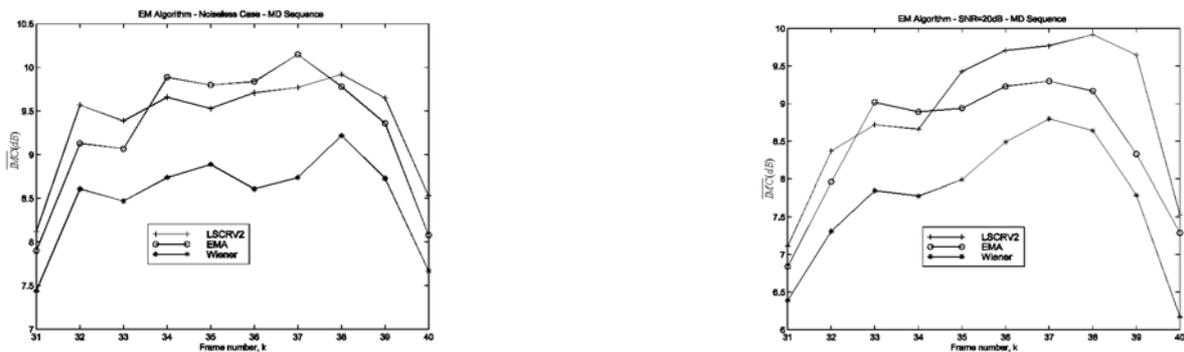


Figure 3. $\overline{IMC}(dB)$ for the noiseless (left) and noisy (right) cases for the “Mother and Daughter” sequence.

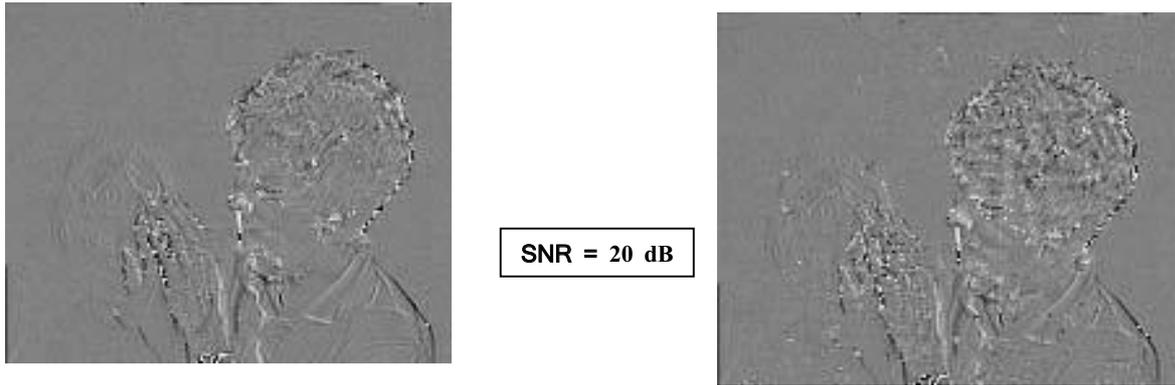


Figure 4. Motion-compensated errors for frame 32 of the “Mother and Daughter” sequence: the Wiener filter (right) and the EMM (left) algorithms.

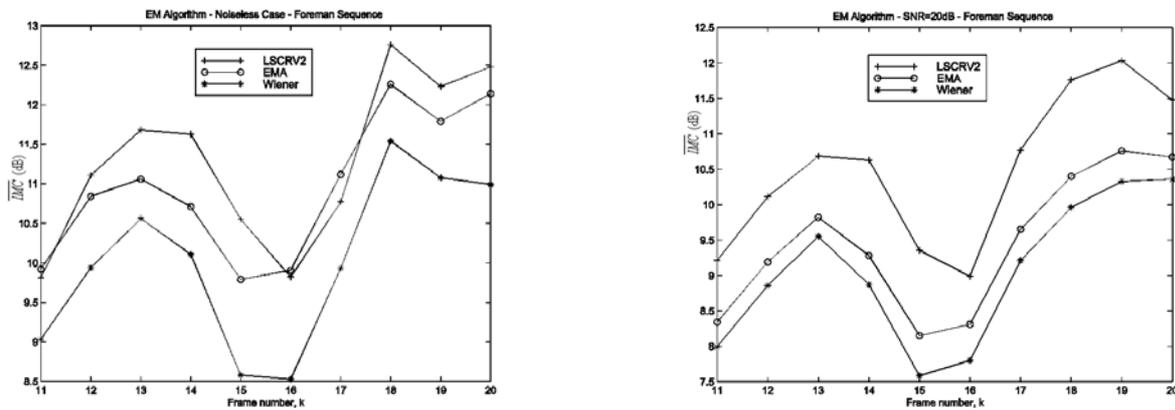


Figure 5. $\overline{MCE}(dB)$ for frames 11-20 of the noiseless (left) and noisy cases with SNR=20dB (right) for the “Foreman” sequence.

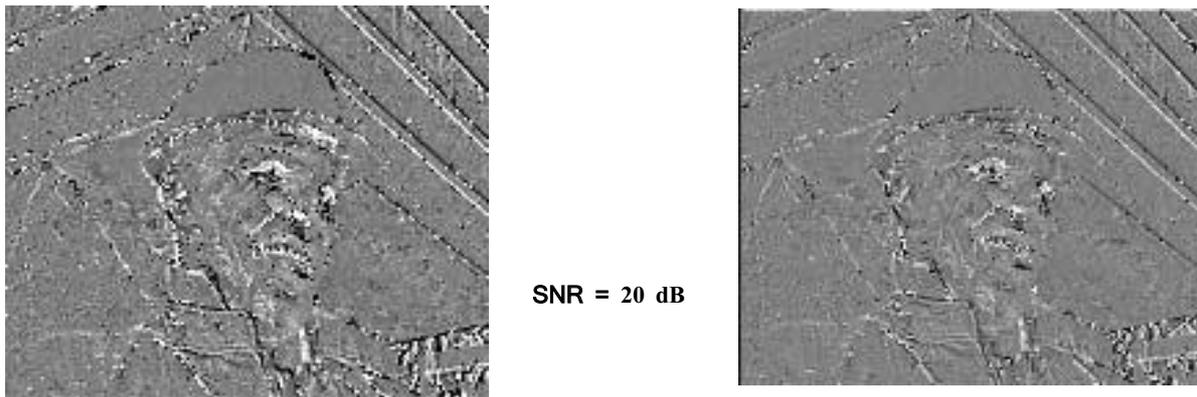


Figure 6. Motion-compensated errors for frame 16 of the “Foreman” sequence: the Wiener filter (left) and the EMM (right) algorithms .