

# Fast Interpolation Operations in Non-Rigid Image Registration

Matthew W. Jacobson<sup>a</sup> and Jeffrey A. Fessler<sup>a\*</sup>

<sup>a</sup>The University of Michigan, EECS Department, Ann Arbor 48109-2122, MI, USA

## ABSTRACT

Much literature on image registration<sup>1-3</sup> has worked with purely geometric image deformation models. For such models, interpolation/resampling operations are often the computationally intensive steps when iteratively minimizing the deformation cost function. This article discusses some techniques for efficiently implementing and accelerating these operations. To simplify presentation, we discuss our ideas in the context of 2D imaging. However, the concepts readily generalize to 3D. Our central technique is a table-lookup scheme that makes somewhat liberal use of RAM, but should not strain the resources of modern processors if certain design parameters are appropriately selected. The technique works by pre-interpolating and tabulating the grid values of the reference image onto a finer grid along one of the axes of the image. The lookup table can be rapidly constructed using FFTs. Our results show that this technique reduces iterative computation by an order of magnitude. When a minimization algorithm employing coordinate block alternation is used, one can obtain still faster computation by storing certain intermediate quantities as state variables. We refer to this technique as *state variable hold-over*. When combined with table-lookup, state variable hold-over reduces CPU time by about a factor two, as compared to table-lookup alone.

## 1. INTRODUCTION AND MOTIVATION

Let  $\mathbb{R}^p$  denote the space of length  $p$  column vectors, for whatever  $p \in \mathbb{N}$ , and similarly let  $\mathbb{R}^{M \times N}$  denote the space of  $M \times N$  real image arrays. This article is concerned with fast methods of deforming and sampling a continuous image  $u(x, y)$ , having the parametric form

$$u(x, y) = \sum_{m=1}^M \sum_{n=1}^N \mathbf{u}_{mn} f(x-m)g(y-n), \quad (1.1)$$

where  $\mathbf{u} \in \mathbb{R}^{M \times N}$  is an array of grid values  $\mathbf{u}_{mn}$ , and  $f: \mathbb{R} \rightarrow \mathbb{R}$  and  $g: \mathbb{R} \rightarrow \mathbb{R}$  are continuous, locally supported interpolating functions, all of which are known. We restrict attention to a parametric family  $\mathcal{W}(\boldsymbol{\alpha})u$  of deformation operations

$$[\mathcal{W}(\boldsymbol{\alpha})u](x, y) \triangleq u(d^X(x, y|\boldsymbol{\alpha}_X), d^Y(x, y|\boldsymbol{\alpha}_Y)). \quad (1.2)$$

Here  $\boldsymbol{\alpha}_X$ ,  $\boldsymbol{\alpha}_Y$ , and  $\boldsymbol{\alpha} = (\boldsymbol{\alpha}_X, \boldsymbol{\alpha}_Y)$  are vectors of deformation parameters and, given the  $x$ -direction deformation basis functions  $\{b_k^X(x, y)\}_{k=1}^{K_x}$  (and similarly for the  $y$ -direction),

$$d^X(x, y|\boldsymbol{\alpha}_X) \triangleq x + \sum_{k=1}^{K_x} \alpha_{X,k} b_k^X(x, y) \quad (1.3)$$

$$d^Y(x, y|\boldsymbol{\alpha}_Y) \triangleq y + \sum_{k=1}^{K_y} \alpha_{Y,k} b_k^Y(x, y). \quad (1.4)$$

Such image transformations are a commonly considered family<sup>2,3</sup> in non-rigid intensity-based image registration. When  $\mathcal{W}(\boldsymbol{\alpha})u$  is sampled at a set of test points  $\{(x_j, y_j)\}_{j=1}^J$  then, combining (1.1) and (1.2), we obtain, for each  $j$ ,

$$\mathcal{W}(\boldsymbol{\alpha})u(x_j, y_j) = \sum_{m,n} f(d^X(x_j, y_j|\boldsymbol{\alpha}_X) - m)g(d^Y(x_j, y_j|\boldsymbol{\alpha}_Y) - n)\mathbf{u}_{mn} \triangleq [\mathbf{W}_{fg}(\boldsymbol{\alpha})\mathbf{u}]_j, \quad (1.5)$$

---

This work was supported in part by NIH/NCI grant 1P01 CA87634.

Thus, deforming and sampling the continuous image  $u$  is equivalent to applying a linear transformation  $\mathbf{W}_{fg}(\boldsymbol{\alpha})\mathbf{u} : \mathbb{R}^{M \times N} \rightarrow \mathbb{R}^J$  to the discrete parameter array  $\mathbf{u}$ . The form of  $\mathbf{W}_{fg}(\boldsymbol{\alpha})$ , however, depends non-linearly on  $\boldsymbol{\alpha}$ . We shall call operators of this form *discrete warping operators*.

Sampling at all  $(x_j, y_j)$  in (1.3) and (1.4) also leads to the quantities

$$\mathbf{d}_j^X(\boldsymbol{\alpha}_X) \triangleq x_j + \sum_{k=1}^{K_x} \alpha_{X,k} b_k^X(x_j, y_j) \quad (1.6)$$

$$\mathbf{d}_j^Y(\boldsymbol{\alpha}_Y) \triangleq y_j + \sum_{k=1}^{K_y} \alpha_{Y,k} b_k^Y(x_j, y_j). \quad (1.7)$$

or in matrix-vector form,

$$\mathbf{d}^X(\boldsymbol{\alpha}_X) \triangleq \mathbf{x} + \mathbf{B}_X \boldsymbol{\alpha}_X \quad (1.8)$$

$$\mathbf{d}^Y(\boldsymbol{\alpha}_Y) \triangleq \mathbf{y} + \mathbf{B}_Y \boldsymbol{\alpha}_Y. \quad (1.9)$$

That is,  $\mathbf{B}_X$  is a  $J \times K_x$  matrix whose columns are the basis function samples  $b_k^X(x_j, y_j)$  (and analogously for  $\mathbf{B}_Y$ ). We refer to  $\mathbf{B}_X$  and  $\mathbf{B}_Y$  as *deformation basis matrices*. We also define  $\mathbf{d}(\boldsymbol{\alpha}) \triangleq [\mathbf{d}^X(\boldsymbol{\alpha}_X) | \mathbf{d}^Y(\boldsymbol{\alpha}_Y)]$  as the  $J \times 2$  matrix whose  $j$ -th row  $\mathbf{d}_j(\boldsymbol{\alpha}) = (\mathbf{d}_j^X(\boldsymbol{\alpha}_X), \mathbf{d}_j^Y(\boldsymbol{\alpha}_Y))$  is the transformation of the test point  $(x_j, y_j)$ .

Based on (1.5), we can see that evaluating operators of the form  $\mathbf{W}_{fg}(\boldsymbol{\alpha})\mathbf{u}$  simply involve interpolating the grid values  $\mathbf{u}_{mn}$  at various locations  $\mathbf{d}_j(\boldsymbol{\alpha})$ . Fast methods for doing so are relevant to registration problems of the form

$$\min. \quad \Phi(\boldsymbol{\alpha}) = \psi(\mathbf{W}_{hh}(\boldsymbol{\alpha})\mathbf{u}) \quad \text{s.t.} \quad \boldsymbol{\alpha} \in \Theta. \quad (1.10)$$

Here  $h$  is an interpolator,  $\Theta$  is a feasible set of deformations, and  $\psi(\cdot) : \mathbb{R}^J \rightarrow \mathbb{R}$  assigns a cost to the deformed image samples  $[\mathbf{W}_{hh}(\boldsymbol{\alpha})\mathbf{u}]_j = \mathcal{W}(\boldsymbol{\alpha})u(x_j, y_j)$  reflecting how well they fit some target data set. Iterative optimization procedures for solving (1.10) will require evaluating  $\Phi$  and some number of its derivatives at various  $\boldsymbol{\alpha}$ . It is clear from (1.10) that it is necessary to carry out the discrete warping operation  $\mathbf{W}_{hh}(\boldsymbol{\alpha})\mathbf{u}$  when evaluating  $\Phi$ . Furthermore, differentiating  $\Phi$  requires additional discrete warpings associated with the derivatives of  $h$ . This is a straightforward consequence of the chain rule. For example, differentiating with respect to  $\boldsymbol{\alpha}_X$  leads to,

$$\begin{aligned} \nabla_{\boldsymbol{\alpha}_X} \Phi(\boldsymbol{\alpha}) &= \nabla_{\boldsymbol{\alpha}_X} \mathbf{d}^X(\boldsymbol{\alpha}_X) \cdot \nabla_{\mathbf{d}^X} \mathbf{W}_{hh}(\boldsymbol{\alpha})\mathbf{u} \cdot \nabla \psi(\mathbf{W}_{hh}(\boldsymbol{\alpha})\mathbf{u}) \\ &= \mathbf{B}_X^T [\mathcal{D} \{ \mathbf{W}_{hh}(\boldsymbol{\alpha})\mathbf{u} \}] \nabla \psi(\mathbf{W}_{hh}(\boldsymbol{\alpha})\mathbf{u}). \end{aligned} \quad (1.11)$$

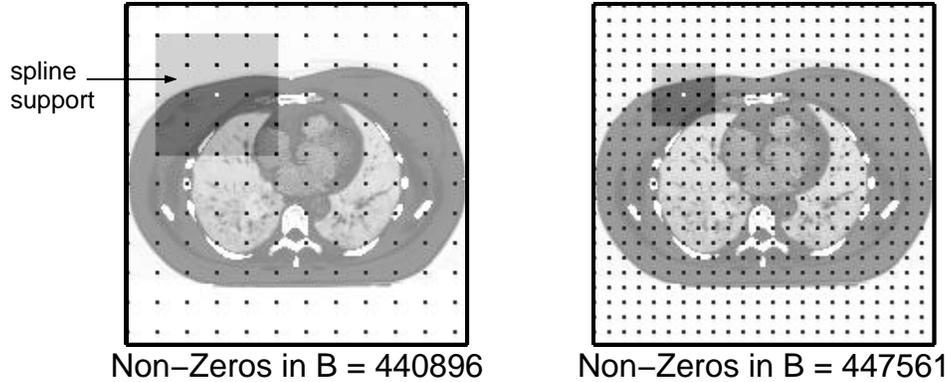
where,  $\mathcal{D} \{z\}$  denotes a the diagonal matrix whose diagonal is the vector  $z$ . Thus, we see that taking a derivative leads to an expression involving not one, but two discrete warping operations,  $\mathbf{W}_{hh}(\boldsymbol{\alpha})\mathbf{u}$  and  $\mathbf{W}_{hh}(\boldsymbol{\alpha})\mathbf{u}$ . Obviously, taking higher derivatives will lead to an even greater preponderance of such operations.

The computational expense of optimization algorithms for (1.10) is therefore dictated by the time needed to manipulate  $\mathbf{B}_X$  (or  $\mathbf{B}_Y$ ), to carry out the various interpolation operations associated with  $\mathbf{W}_{hh}(\boldsymbol{\alpha})\mathbf{u}$  and its derivatives, and to evaluate the requisite derivatives of  $\psi(\cdot)$ . As we shall discuss, the manipulation of  $\mathbf{B}_X$  is typically a minor computational issue compared to the discrete warpings. These matrices usually have a small number of non-sparse elements or an otherwise favorable structure.

The computational expense of manipulating  $\psi(\cdot)$  is problem-specific. In this article, we shall be concerned mainly with the least squares cost function

$$\Phi_{LS}(\boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{W}_{hh}(\boldsymbol{\alpha})\mathbf{u} - \mathbf{v}\|^2. \quad (1.12)$$

That is,  $\psi(\mathbf{z}) = \|\mathbf{z} - \mathbf{v}\|^2/2$ , where  $\|\cdot\|$  denotes the standard Euclidean norm and  $\mathbf{v} \in \mathbb{R}^J$  is a given vector of target image values. Techniques for minimizing (1.12) have been explored in various works.<sup>1-3</sup> For this problem, the form of  $\psi(\cdot)$  and its derivatives is trivial compared to the discrete warpings. Thus, the discrete warpings are the dominant computational burden, and the acceleration techniques that we propose make a significant difference in overall CPU time. For other  $\psi(\cdot)$ , the acceleration techniques may also be helpful to some degree.



**Figure 1.** Depiction of a typical system of deformation basis functions, both coarse and fine, and the relationship of its geometry to the grid of  $\mathbf{u}$ . In this illustration,  $\mathbf{u}$  is a  $184 \times 184$  torso phantom. Irrespective of the fineness of the deformation grid, the number of non-zero entries in  $\mathbf{B}_C$ ,  $C \in \{X, Y\}$  is about 450000.

In the next section, we discuss, in a bit more detail, the general considerations affecting the speed of the discrete warping operations. In Section 3, we present some techniques for accelerating the interpolation operations involved. One of these techniques, which we call *state variable hold-over*, applies to algorithms that alternately update  $\alpha_X$  and  $\alpha_Y$ , and is based on saving in memory certain intermediate quantities associated with whichever variable is being held fixed. The second technique is a table-lookup scheme in which  $\mathbf{u}$  is pre-interpolated at small intervals along (at least) one axis. As we show, the pre-interpolation can be done rapidly using FFT operations. Finally, in Section 4, we test our various techniques on the registration of two simulated torso phantoms.

## 2. GENERAL COMPUTATIONAL ASPECTS

### 2.1. The Role of the Deformation Basis Matrix

It is apparent from (1.11), that differentiating the cost function  $\Phi$  requires multiplication by  $\mathbf{B}_X^T$  and  $\mathbf{B}_Y^T$ . In addition, it is important to remember that a discrete warping  $\mathbf{W}_{fg}(\alpha)\mathbf{u}$  depends on  $\alpha$  only through  $d^X(\alpha_X)$  and  $d^Y(\alpha_Y)$  (see (1.5)), which in turn depend on  $\alpha$  only through  $\mathbf{B}_X\alpha_X$  and  $\mathbf{B}_Y\alpha_Y$  (see (1.8) and (1.9)), respectively. It is therefore pertinent to consider the structure of  $\mathbf{B}_X$  and  $\mathbf{B}_Y$  and the computational demands that manipulating them entails.

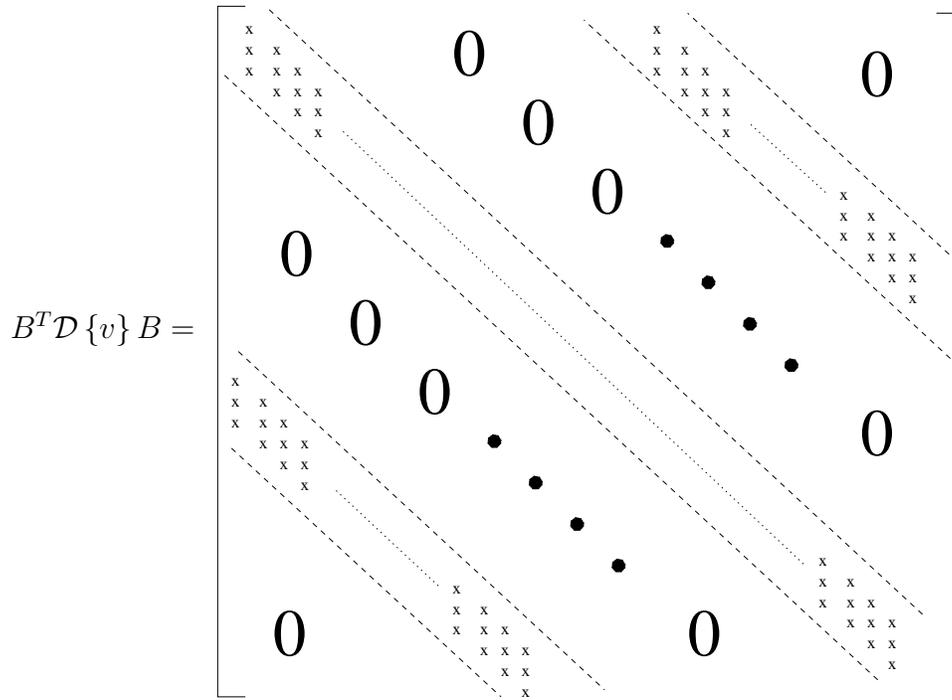
At this time, it is convenient to define

$$\mathbf{B} = \begin{bmatrix} \mathbf{B}_X & 0 \\ 0 & \mathbf{B}_Y \end{bmatrix}.$$

As mentioned, the columns of  $\mathbf{B}_X$  are the samples  $b_k^X(x_j, y_j)$  of the  $x$ -direction basis functions (and similarly for  $\mathbf{B}_Y$ ). The form of  $\mathbf{B}$  is therefore dictated by the choice of deformation basis functions. A typical scheme by which basis functions are chosen is depicted in Figure 1. Shifted copies of locally supported functions (typically tensor products of B-splines) are centered at the nodes of a deformation space grid, usually a coarse subset of the grid of the  $\mathbf{u}_{mn}$  parameters. One makes the grid finer (reducing the support of the basis functions accordingly) to span finer deformations. Since the entries of  $\mathbf{B}_X$  and  $\mathbf{B}_Y$  are derived from shifted-copies of functions, these matrices have a circulant-like structure.

When the test points  $\{(x_j, y_j)\}_{j=1}^J$  are simply the set of integer coordinates  $\{(m, n)\}$ ,  $1 \leq m \leq M$ ,  $1 \leq n \leq N$ , the number of non-zero entries in  $\mathbf{B}_X$  (or  $\mathbf{B}_Y$ ) is approximately  $MN$  times the area (as measured in number of nodes) of each of the basis functions. (Note that this number is constant irrespective of the fineness of the deformation space grid). Thus, for example, when the basis functions are cubic B-splines, the number of non-zero entries of  $\mathbf{B}_X$  is about  $16MN$  and the matrix can easily be stored in RAM. In 3D problems, the storage requirements may be prohibitive. However, because the  $\mathbf{B}_X$  and  $\mathbf{B}_Y$  matrices have a circulant-like structure, fast convolution techniques can be used to multiply with these matrices and their transposes.

When a fine deformation grid is used, the support of the basis functions is relatively small. Together with the circulant-like structure of  $\mathbf{B}_X$  and  $\mathbf{B}_Y$ , this implies that matrix products of the form  $\mathbf{B}_C^T \mathbf{D}\{z\} \mathbf{B}_C$ ,  $C \in \{X, Y\}$ , and hence also



**Figure 2.** Depiction of the multi-banded structure of  $B^T \mathcal{D}\{z\} B$ .

those of the form  $B^T \mathcal{D}\{z\} B$ , have a sparse, multi-banded structure as depicted in Figure 2. Such matrices are rapidly inverted by sparse implementations of Gaussian-elimination. The advantage of this in the least squares problem will be made apparent in Section 2.3.

Manageable  $B$  matrices are also obtained when locally supported basis functions (of whatever form) are situated at only a few select locations in the image (such as when the significant deformation is known to occur in a particular region). They would be manageable simply because their size and number of non-zero entries may be comparatively small.

## 2.2. The Constituents of Interpolation/Sampling Operations

As mentioned, the evaluation of a discrete warping operator  $\mathbf{W}_{fg}(\alpha)\mathbf{u}$  involves first computing  $d^X(\alpha_X)$  and  $d^Y(\alpha_Y)$  and then (see (1.5)) interpolating the parameters  $\mathbf{u}_{mn}$  at every deformed test point  $\mathbf{d}_j(\alpha) = (d_j^X(\alpha_X), d_j^Y(\alpha_Y))$ . Some of the ideas relevant to doing so are depicted in Figure 3.

There are four principal steps to each interpolation. Firstly, one locates the region of non-trivial grid values  $\mathbf{u}_{mn}$  that contribute to the interpolated value  $[\mathbf{W}_{fg}(\alpha)\mathbf{u}]_j$ . Because multiplicatively separable interpolation is assumed here and because  $f$  and  $g$  are locally supported interpolating functions, there will be some local rectangular region of grid values  $\mathbf{u}_{mn}$  near each given  $\mathbf{d}_j$  that contribute to  $[\mathbf{W}_{fg}(\alpha)\mathbf{u}]_j$ . This region will be the intersection of a  $\alpha_X$  dependent strip and a  $\alpha_Y$  dependent strip as shown in Figure 3. The width of the strip is the width of  $f$  or  $g$ , as appropriate. Henceforth, we shall call such a region an *interpolation rectangle*. Next, one accesses the associated grid values  $\mathbf{u}_{mn}$  in this region from memory. As this is done for every  $j$ , many non-contiguous memory accesses are required. Therefore, this step incurs appreciable CPU time. Thirdly, the interpolation weights applied to the different  $\mathbf{u}_{mn}$  is computed. Typically, this involves the repeated evaluation of non-linear or transcendental functions. Therefore, this step also incurs appreciable CPU time. Finally, the weights are applied to the  $\mathbf{u}_{mn}$  and summed to give  $[\mathbf{W}_{fg}(\alpha)\mathbf{u}]_j$ . The computational effort in all of these steps is proportional to the number of grid points  $\mathbf{u}_{mn}$  in each interpolation rectangle.

## 2.3. The Case of the Least-Squares Problem

In this section, we look at the form of the derivatives of the least squares cost function  $\Phi_{LS}$  and see more concretely why, for this problem, discrete warpings are the computational bottleneck in derivative-based minimization methods. We assume

throughout that the interpolator  $h$  is as many times differentiable as the context requires.

Any first order method (e.g., steepest descent) will rely on the gradient, which has the form

$$\nabla \Phi_{\text{LS}}(\boldsymbol{\alpha}) = \mathbf{B}^T \begin{bmatrix} \mathbf{g}_X(\boldsymbol{\alpha}) \\ \mathbf{g}_Y(\boldsymbol{\alpha}) \end{bmatrix}, \quad (2.1)$$

where,

$$\begin{aligned} \mathbf{e}(\boldsymbol{\alpha}) &= \mathbf{W}_{hh}(\boldsymbol{\alpha})\mathbf{u} - \mathbf{v} \\ \mathbf{g}_X(\boldsymbol{\alpha}) &= \mathcal{D}\{\mathbf{e}(\boldsymbol{\alpha})\} \mathbf{W}_{hh}(\boldsymbol{\alpha})\mathbf{u} \\ \mathbf{g}_Y(\boldsymbol{\alpha}) &= \mathcal{D}\{\mathbf{e}(\boldsymbol{\alpha})\} \mathbf{W}_{hh}(\boldsymbol{\alpha})\mathbf{u}. \end{aligned}$$

Furthermore, any second order (e.g., Newton's) method will require not only the gradient, but also the Hessian, given by

$$\nabla^2 \Phi_{\text{LS}}(\boldsymbol{\alpha}) = \mathbf{B}^T \begin{bmatrix} \mathbf{H}_{XX}(\boldsymbol{\alpha}) & \mathbf{H}_{XY}(\boldsymbol{\alpha}) \\ \mathbf{H}_{XY}^T(\boldsymbol{\alpha}) & \mathbf{H}_{YY}(\boldsymbol{\alpha}) \end{bmatrix} \mathbf{B}, \quad (2.2)$$

where,

$$\begin{aligned} \mathbf{H}_{XX} &= \mathcal{D}\{\mathbf{W}_{hh}\mathbf{u}\} \mathcal{D}\{\mathbf{e}\} + \mathcal{D}\{\mathbf{W}_{hh}\mathbf{u}\}^2 \\ \mathbf{H}_{YY} &= \mathcal{D}\{\mathbf{W}_{hh}\mathbf{u}\} \mathcal{D}\{\mathbf{e}\} + \mathcal{D}\{\mathbf{W}_{hh}\mathbf{u}\}^2 \\ \mathbf{H}_{XY} &= \mathcal{D}\{\mathbf{W}_{hh}\mathbf{u}\} \mathcal{D}\{\mathbf{e}\} + \mathcal{D}\{\mathbf{W}_{hh}\mathbf{u}\} \mathcal{D}\{\mathbf{W}_{hh}\mathbf{u}\}. \end{aligned}$$

We see that the expressions for the various derivatives are replete with discrete warping operations. Moreover, the only other operations of any significance are the multiplications with  $\mathbf{B}$  and its constituents. As mentioned, this was to be expected based on the simple form of  $\psi(\cdot)$  for the problem.

What is noteworthy beyond this, however, is that the the constituents of the Hessian  $\mathbf{B}_X^T \mathbf{H}_{XX} \mathbf{B}_X$ ,  $\mathbf{B}_Y^T \mathbf{H}_{YY} \mathbf{B}_Y$ , ... have the multi-banded form of Figure 2, since  $\mathbf{H}_{XX}$ ,  $\mathbf{H}_{YY}$ , and  $\mathbf{H}_{XY}$  are all diagonal. Consequently, the total Hessian  $\nabla^2 \Phi_{\text{LS}}$  also has this structure. This means that, not only do the manipulations of  $\mathbf{B}$  not contribute significantly to computational overhead, but Hessian inversions, when used in Newton-like methods, should also not contribute significantly. The inversions may be done rapidly with Gaussian elimination.

### 3. ACCELERATION TECHNIQUES

#### 3.1. State Variable Hold-Over

In this section, we discuss an acceleration technique that we call *state variable hold-over*. The technique is applicable to algorithms that alternately update  $\boldsymbol{\alpha}_X$  and  $\boldsymbol{\alpha}_Y$  as follows:

ALGORITHM 3.1. [*Skeleton of a block alternating algorithm*]

*Initialization.* Choose a number of iterations  $L > 0$  and sub-iteration parameters  $L_X, L_Y > 0$ .

For  $i = 0, \dots, L$

For  $k := 0, \dots, L_X - 1$

Update  $\boldsymbol{\alpha}_X$  only:  $(\boldsymbol{\alpha}_X^{\text{old}}, \boldsymbol{\alpha}_Y^{\text{old}}) \mapsto (\boldsymbol{\alpha}_X^{\text{new}}, \boldsymbol{\alpha}_Y^{\text{old}})$

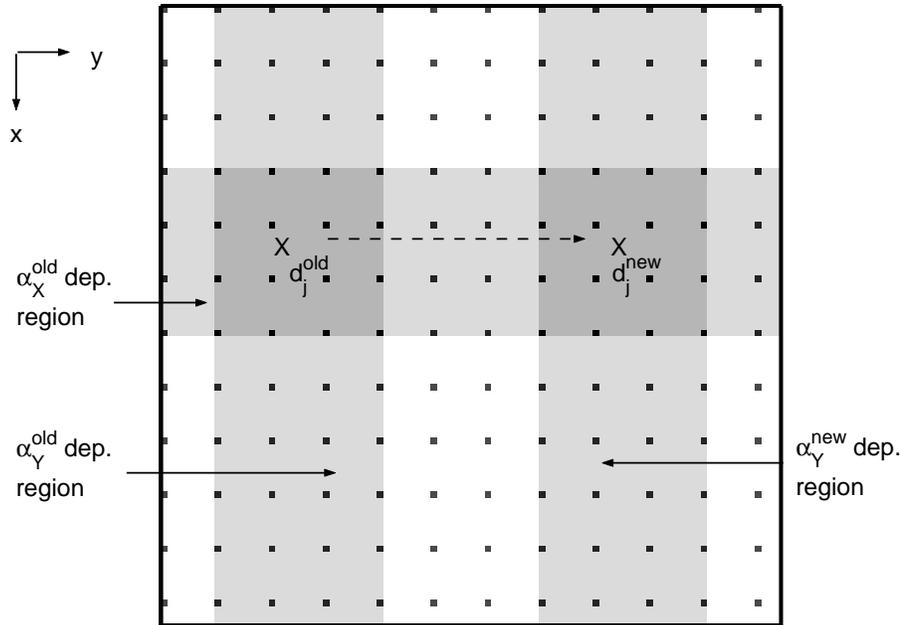
end

For  $k := 0, \dots, L_Y - 1$

Update  $\boldsymbol{\alpha}_Y$  only:  $(\boldsymbol{\alpha}_X^{\text{old}}, \boldsymbol{\alpha}_Y^{\text{old}}) \mapsto (\boldsymbol{\alpha}_X^{\text{old}}, \boldsymbol{\alpha}_Y^{\text{new}})$

end

end



**Figure 3.** Locating relevant grid neighborhoods in interpolation operations relevant to discrete warpings. Also, a depiction of how state variable hold-over facilitates this.

When certain variables are held fixed, quantities that are functions of the fixed variables can be held in memory and re-used each time the non-fixed variables are updated. Thus, additional, consecutive updates of the non-fixed variable are more cheaply obtained than if all quantities were freshly computed. How this concept applies to the interpolation operations presently discussed is illustrated in Figure 3 for the case where  $\alpha_X$  is the fixed variable. There, we see that, when  $\alpha_Y$  is updated, the  $x$ -coordinate of  $d_j^{new}$  is the same as of  $d_j^{old}$ . Thus, the local interpolation rectangle has the same  $x$ -coordinates as for  $d_j^{old}$ . Locating the new interpolation rectangle only requires computing the new  $y$ -coordinates. Similarly, the  $x$ -direction interpolation weights applied to the new rectangle is the same as for the old. Hence, only the  $y$ -direction weights need to be recomputed.

### 3.2. Table-Lookup Elimination

In this section, we propose a table-lookup technique that eliminates the need to interpolate along at least one axis. The idea is to pre-interpolate along that axis at a small sampling interval, compared to the original grid spacing of  $u$ , and tabulate the results. During the iterations of the minimization algorithm, interpolation along this axis may then be replaced with nearest-neighbor table-lookup. Figure 4 illustrates this for when the  $x$ -axis (but not the  $y$ -axis) is finely sampled. As shown in the figure, the number of grid values in each local interpolation rectangle reduces as compared to Figure 3 when no pre-interpolation is done. The reduction factor is the support width of the interpolator for the finely sampled direction.<sup>†</sup> Accordingly, the overall time to compute  $W_{fg}(\alpha)u$  reduces by the same factor.

#### 3.2.1. Incorporating State Variable Hold-Over

A further advantage to the table-lookup scheme becomes apparent when one considers incorporating state variable hold-over (see Section 3.1). Suppose, as in Figure 4, that  $\alpha_X$ , but not  $\alpha_Y$ , was being updated as in Algorithm 3.1. If state variable hold-over is used, the weights for the  $y$ -direction interpolation (which is the only direction in which interpolation occurs) are held in memory throughout the consecutive updates of  $\alpha_X$ . Thus, updating  $W_{fg}(\alpha)u$  simply involves multiplying the values in the new interpolation rectangle by these already available weights and summing. No new weight calculations are required.

<sup>†</sup>The effect of this reduction will be particularly prominent if the interpolator for the finely sampled direction is very wide compared to those of other directions.

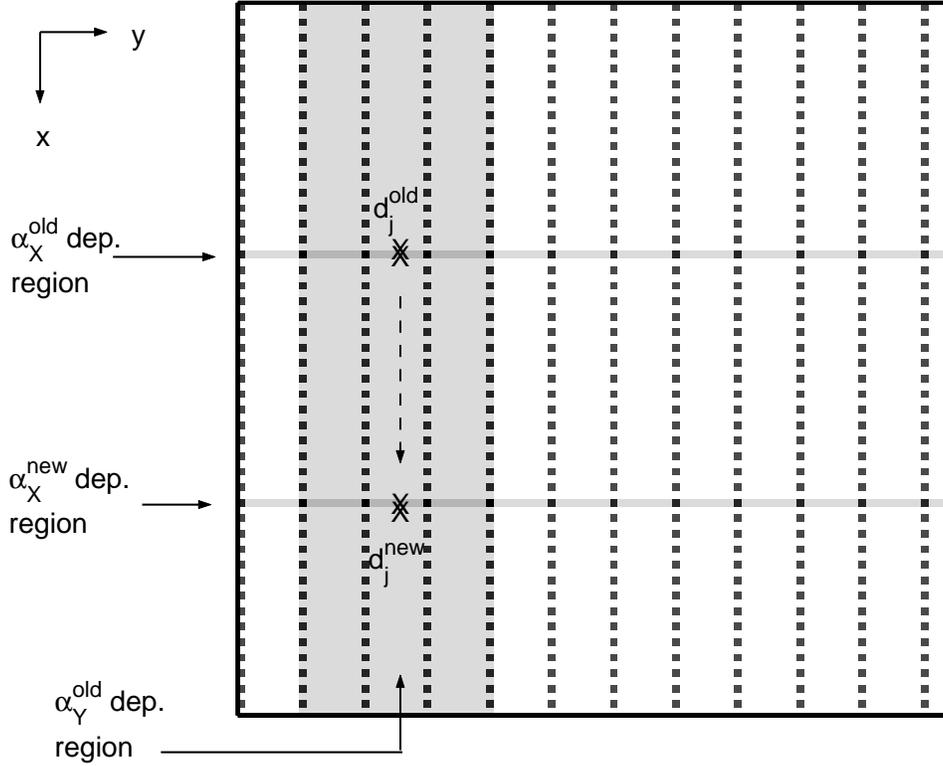


Figure 4. Illustration of table-lookup elimination.

### 3.2.2. Generating the Table

The table-assisted interpolations in Figure 4 approximate each  $[\mathbf{W}_{fg}(\boldsymbol{\alpha})\mathbf{u}]_j$  as follows (cf. Equation (1.5)),

$$[\mathbf{W}_{fg}(\boldsymbol{\alpha})\mathbf{u}]_j \approx \sum_{n=1}^N g(y-n) \left[ \sum_{p=1}^M f(x-p)\mathbf{u}_{pn} \right] \Bigg|_{\substack{x=\langle \mathbf{d}_j^X(\boldsymbol{\alpha}_X) \rangle \\ y=\mathbf{d}_j^Y(\boldsymbol{\alpha}_Y)}} \quad (3.1)$$

Here  $\langle \mathbf{d}_j^X \rangle$  is  $\mathbf{d}_j^X$  rounded to the nearest in a set of finely spaced  $x$ -values  $\{x_{ms}\}$  where

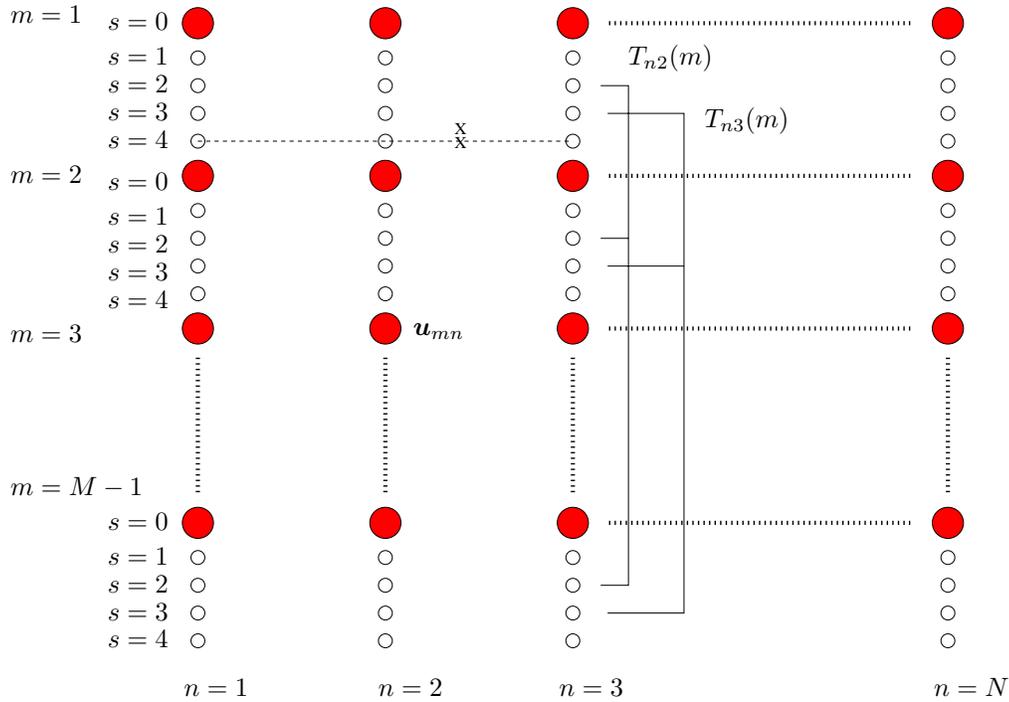
$$x_{ms} = m + s/S, \quad 1 \leq m \leq M-1, 0 \leq s \leq S-1. \quad (3.2)$$

The parameter  $S \in \mathbb{N}$  is a pre-selected value that controls the fineness of the sampling. This indexing scheme is illustrated in Figure 5.

The table-lookup technique assumes that the quantity  $\sum_{p=1}^M f(x-p)\mathbf{u}_{pn}$  has been pre-tabulated at all  $x_{ms}$ . A naive way of constructing such a table would be to loop over the values in  $\{x_{ms}\}$  in ascending order and, for each one, to compute the table entries

$$T_{mns} \triangleq \sum_{p=1}^M f(x_{ms}-p)\mathbf{u}_{pn} \quad (3.3)$$

directly. It proves more efficient, however, to process  $\{x_{ms}\}$  in integer-spaced groups, each corresponding to a fixed  $s$



**Figure 5.** By taking integer-spaced sets of samples, the lookup table can be rapidly constructed via FFTs.

value. Defining  $\kappa_{f,s}(m) \triangleq f(m + s/S)$  and incorporating (3.2) into (3.3),

$$\begin{aligned}
 T_{ns}(m) &= \sum_p f(m + s/S - p) \mathbf{u}_{pn} \\
 &= \sum_p \kappa_{f,s}(m - p) \mathbf{u}_n(p) \\
 &= [\kappa_{f,s} \otimes_p \mathbf{u}_n](m)
 \end{aligned} \tag{3.4}$$

where  $\otimes_p$  denotes 1D discrete convolution with respect to  $p$ . The discrete convolution in (3.4) can be done rapidly via FFTs. Moreover, the FFT of each  $\kappa_{f,s}$  depends only on  $f$  and can be archived in permanent storage.

The choice of  $S$  reflects a compromise between memory consumption and accuracy. Obviously, for larger  $S$ , the spacing of the  $\{x_{ms}\}$  is smaller and the approximation accuracy of (3.1) increases. However, the storage requirements of the implementation is also magnified by  $S$ . The total number of stored grid parameters  $\mathbf{u}_{mn}$  effectively increases to  $SMN$  and the number of data in the FFT of the  $\kappa_{f,s}$  is  $O(SM)$ . Since  $f$  and  $g$  are typically gradually varying functions, however, a moderate  $S$  value ( $\leq 10$ ) is likely to be sufficient.

#### 4. COMPUTATIONAL TESTS

We applied the acceleration techniques of Section 3 to the registration of the  $128 \times 128$  torso image shown in Figure 6(a) with the image shown in Figure 6(b). The absolute difference between these images appears in Figure 6(d). The array of pixel values in Figure 6(a) played the role of the reference image data,  $\mathbf{u}$ . By applying a certain warping  $\mathbf{W}_{hh}(\alpha^{\text{true}}) \mathbf{u}$  to  $\mathbf{u}$ , we generated the target image data vector  $\mathbf{v} \in \mathbb{R}^{128^2}$  which, when rearranged as a  $128 \times 128$  image, gives the image displayed in Figure 6(b). The test points  $\{(x_j, y_j)\}_{j=1}^J$ ,  $J = 128^2$ , coincided with the grid points  $\{(m, n) : 1 \leq m \leq 128, 1 \leq n \leq 128\}$ . The interpolator  $h$  was the cubic B-spline functions, each 4 pixels wide. For deformation basis functions, we used 2D tensor products of cubic B-splines, also of width 4. The basis functions  $\{b_k^X(x, y)\}_{k=1}^{120}$  were each centered at a pixel in the shaded region about the left lung boundary in Figure 6(c). Similarly, deformation basis functions  $\{b_k^Y(x, y)\}_{k=1}^{283}$  were situated at pixels in the shaded region by the right lung. The true parameter values  $\alpha_{X,k}^{\text{true}}$  and  $\alpha_{Y,k}^{\text{true}}$  were all equal to -10 pixel lengths.

To minimize  $\Phi_{LS}$ , we employed the following special case of Algorithm 3.1:

ALGORITHM 4.1. [A Block-Alternating Newton-like Method]

For  $i = 0, \dots, L$

For  $C \in \{X, Y\}$

For  $k = 0, \dots, L_C - 1$

$$\begin{aligned}\mathcal{H} &:= \mathbf{B}_C^T \mathbf{H}_{CC}(\alpha_X, \alpha_Y) \mathbf{B}_C \\ \mathbf{g} &:= \mathbf{B}_C^T \mathbf{g}_C(\alpha_X, \alpha_Y) \\ \alpha_C &:= \alpha_C - \{\mathcal{H}\}_+^{-1} \mathbf{g}\end{aligned}\tag{4.1}$$

end

end

end

In (4.1), the transformation  $\{\cdot\}_+$  is some operation that makes  $\mathcal{H}$  positive definite. When  $\{\cdot\}_+$  is simply the identity mapping,  $\{\mathcal{H}\}_+ = \mathcal{H}$ , these update equations reduce to Newton steps in  $\alpha_X$  and  $\alpha_Y$  respectively. However,  $\Phi_{LS}$  is non-convex, so  $\{\cdot\}_+$  would normally be something non-trivial.

In the present experiment,  $\{\mathcal{H}\}_+$  was the transformation that alters the eigendecomposition of  $\mathcal{H}$  by replacing each eigenvalue,  $\lambda_k$ , with  $\max\{|\lambda_k|, 0.1\}$ .<sup>‡</sup> This can be viewed as a refinement of the Levenberg-Marquardt method.<sup>§</sup> In the Levenberg-Marquardt method, the Hessian eigenvalues are all increased by a common constant. Conversely, in our method, eigenvalues that are sufficiently large to begin with are left alone, thus avoiding unnecessarily large curvatures and slower convergence. The computation of  $\{\mathcal{H}\}_+^{-1}$  is significantly more time-consuming than for the  $\{\cdot\}_+$  transformation used in the Levenberg-Marquardt method. However, another advantage to the block alternation technique is that the dimensions of  $\mathcal{H}$  are smaller than the full Hessian  $\nabla^2 \Phi_{LS}(\alpha)$ , making the transformation more manageable.

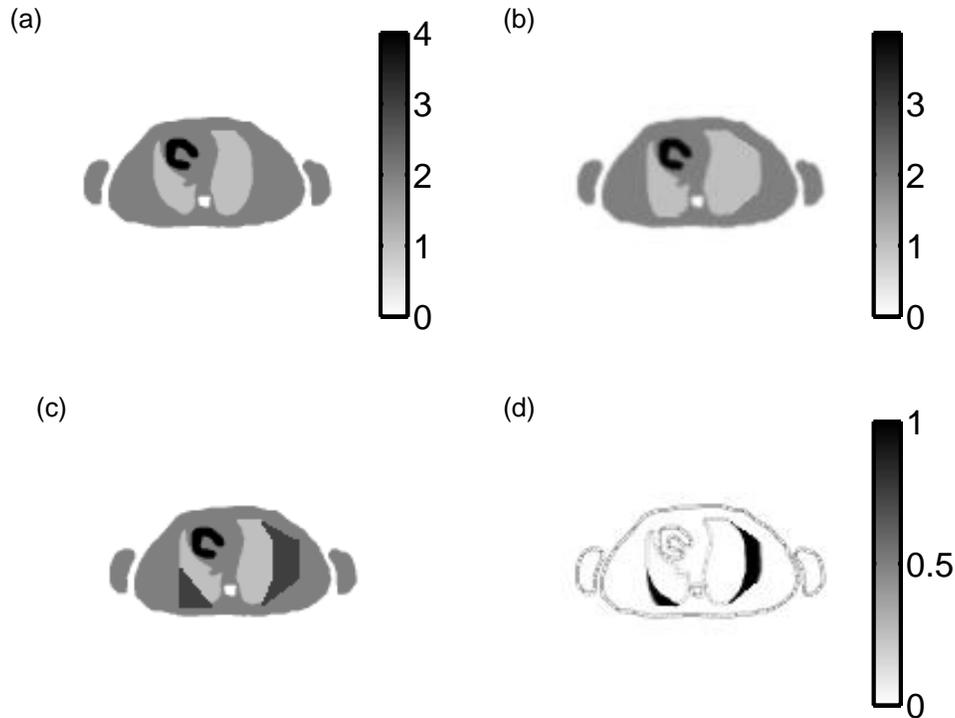
We implemented Algorithm 4.1 using all combinations of the acceleration techniques in Section 3. Each implementation was then tested on the reference and target image data. In all cases, the iteration/sub-iteration loop sizes were  $L = 5$  and  $L_X = L_Y = 10$ . Also, the components of the initial  $\alpha$  were set, in each test, to -5 pixel lengths. The implementations employing the table-lookup technique did so with  $S = 10$ . Two sets of tables were maintained. One set was derived from pre-interpolation (with  $h$  and all relevant derivatives) in the  $x$ -direction. The second set was similarly derived by pre-interpolating in the  $y$ -direction.<sup>¶</sup>

In Figures 7 and 8, the progress of each algorithm implementation is plotted versus CPU time. Figure 7 only accounts for the time spent on gradient and Hessian evaluations, which is where all of the interpolation operations occur. Figure 8 plots descent against the total CPU time expended up to that point in the program execution. The latter took into account overhead for generating the lookup tables before the start of the iterations and computing  $\{\mathcal{H}\}_+^{-1}$ . The trials marked ‘Plain’ refer to the case where none of the acceleration techniques were used, i.e., interpolations were done entirely-on-the-fly. The memory consumption, measured in millions of stored numbers, was also determined for each implementation (see Table 1). For 3D registration, one would require this amount of memory for each image slice.

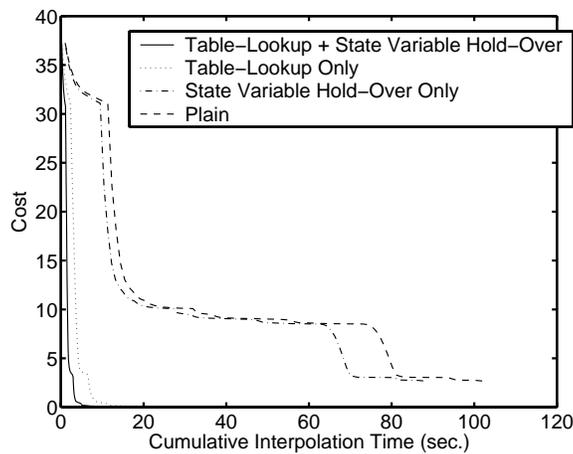
<sup>‡</sup>The threshold value of 0.1 was empirically chosen. More generally, one desires a constant threshold that is small compared to the Hessian eigenvalues in the neighbourhood of the solution.

<sup>§</sup>The application of Levenberg-Marquardt, without coordinate block alternation, was considered for the least squares registration problem.<sup>2,3</sup>

<sup>¶</sup>Clearly, many variants of this scheme are possible.



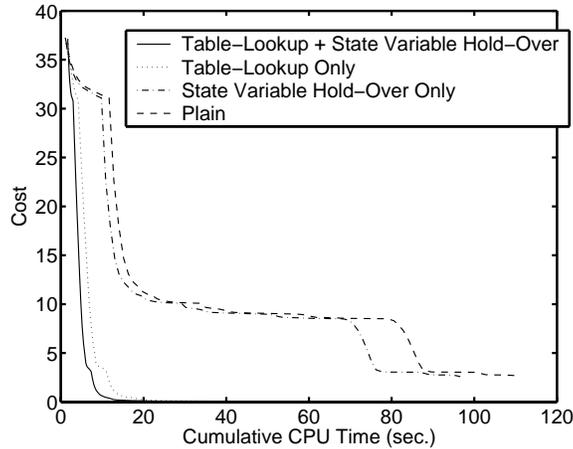
**Figure 6.** (a) The reference image. (b) The target image. (c) Locations of the deformation basis functions. Cubic B-splines were placed at all pixels in the darkened regions near the lung boundaries. Basis functions  $\{b_k^X\}_{k=1}^{K_x}$  were placed in the region at the left lung and  $\{b_k^Y\}_{k=1}^{K_y}$  were placed in the region at the right lung. (d) Absolute difference between the reference and target images.



**Figure 7.** The progress of the algorithm implementations versus cumulative time spent on interpolation operations, i.e., on gradient and Hessian computations.

## 5. CONCLUSIONS

From Figures 7 and 8, we see that the table-lookup strategy decreases the computation time, as predicted, by a factor of 4, the width of the B-spline interpolators. When table-lookup is used, adding the technique of state variable hold-over leads to additional acceleration by about a factor of 2. For the implementations that used table lookup, a delayed start of less than 2 seconds is seen in Figure 8. This delay is due to the time spent generating tables, but constitutes a small fraction of the overall iteration time.



**Figure 8.** The progress of the algorithm implementations versus total CPU time.

**Table 1.** Memory consumption for the different implementations.

Method	Consumed RAM (millions of numbers/slice)
Table-Lookup + State Variable Hold-Over	1.65
Table-Lookup Only	1.38
State Variable Hold-Over Only	0.89
Plain	0.28

We can expect still faster computation per iteration when using the Levenberg-Marquardt method to perform the updates in (4.1). This is because Levenberg-Marquardt updates rely only on Gaussian elimination, and will be rapid since the transformation  $\{\cdot\}_+$  employed by Levenberg-Marquardt preserves the banded structure of the partial Hessian  $\mathcal{H}$  (cf. Section 2.1). However, the descent may be slower with Levenberg-Marquardt than with the  $\{\cdot\}_+$  that we have used here.

The acceleration comes at the expense of extra memory consumption. With  $128 \times 128$  slices and  $S = 10$ , Table 1 shows that the memory requirement by no means strains the resources of modern processors. For larger data sets  $\mathbf{u}$ , it may be necessary to use smaller  $S$ . A variety of other memory-conserving compromises are also possible. For example, the tabulated data could be approximated as 2 byte short integers, as opposed to floating point. Or, instead of maintaining a separate table of finely spaced 1D grid data (as in Figure 4) for each Cartesian direction, one could do so for only one direction. In our future work, we plan to test the impact of these various compromises. Conversely, however, if memory resources are abundant and  $\mathbf{u}$  is not too large, one could contemplate even more liberal use of memory than was done here. One might even generate a table that finely samples  $\mathbf{u}$  in all dimensions at once (unlike in Figure 4). In that case, *all* interpolation operations would be eliminated, approximated instead by nearest-neighbor table-lookup.

### Acknowledgements

The first author wishes to thank Michael Zibulevsky for inspiring the table lookup ideas discussed in this article.

### REFERENCES

1. M. W. Jacobson and J. A. Fessler, "Joint estimation of image and deformation parameters in motion-corrected PET," in *Proc. IEEE Nuc. Sci. Symp. Med. Im. Conf.*, **5**, pp. 3290–4, 2003.
2. J. Kybic, P. Thevenaz, A. Nirkko, and M. Unser, "Unwarping of unidirectionally distorted EPI images," *IEEE Tr. Med. Imag.* **19**, pp. 80–93, Feb. 2000.
3. J. Kybic and M. Unser, "Fast parametric elastic image registration," *IEEE Tr. Im. Proc.* **12**, pp. 1427–42, Nov. 2003.