# Non-negative Tensor Patch Dictionary Approaches for Image Compression and Deblurring Applications[*]

Elizabeth Newman[†] and Misha E. Kilmer[‡]

**Abstract.** In recent work (Soltani, Kilmer, Hansen, BIT 2016), an algorithm for non-negative tensor patch dictionary learning in the context of X-ray CT imaging and based on a tensor-tensor product called the $t$-product (Kilmer and Martin, 2011) was presented. Building on that work, in this paper, we use of non-negative tensor patch-based dictionaries trained on other data, such as facial image data, for the purposes of either compression or image deblurring. We begin with an analysis in which we address issues such as suitability of the tensor-based approach relative to a matrix-based approach, dictionary size and patch size to balance computational efficiency and qualitative representations. Next, we develop an algorithm that is capable of recovering non-negative tensor coefficients given a non-negative tensor dictionary. The algorithm is based on a variant of the Modified Residual Norm Steepest Descent method. We show how to augment the algorithm to enforce sparsity in the tensor coefficients, and note that the approach has broader applicability since it can be applied to the matrix case as well. We illustrate the surprising result that dictionaries trained on image data from one class can be successfully used to represent and compress image data from different classes and across different resolutions. Finally, we address the use of non-negative tensor dictionaries in image deblurring. We show that tensor treatment of the deblurring problem coupled with non-negative tensor patch dictionaries can give superior restorations as compared to standard treatment of the non-negativity constrained deblurring problem.

**Key words.** tensor, patch dictionary, image compression, image deblurring, MRNSD, sparsity constraint

**AMS subject classifications.** 65F22, 65F99, 65N20, 65N21

**1. Introduction.** Many applications in imaging science, such as image deblurring and image reconstruction, typically model the object to be recovered as a vector of unknowns, and the forward operator as a matrix. Treating image and video data processing problems using tensor approaches is yet a new, but increasingly popular and promising approach, as suggested by recent literature [4, 13, 12, 5, 23, 19, 17, 22, 24]. However, a close look at the increasing body of literature in which tensor decompositions are used in practice shows that no one specific tensor decomposition has fit all these image and video applications equally well. Indeed, here, as in many other multiway data processing and representation problems, the type of decomposition to be employed may be quite specific to the application.

Decompositions based on a tensor-tensor product called the $t$-product [9] have proven to

---

   [†]Department of Mathematics, Emory University, Atlanta, GA (elizabeth.newman@emory.edu)

   [‡]Department of Mathematics, Tufts University, Medford, MA (misha.kilmer@tufts.edu, http://www.tufts.edu/~mkilme01/).

be particularly useful in applications where there is a natural orientation dependence to be preserved, such as pixel or voxel position, relative to some other variable such as number of images or time (see [4, 24] for example). The $t$-product is advantageous over other tensor decompositions because of the algebraic framework induced by the definition of the $t$-product, which enables the definition and computation of factorizations reminiscent of their matrix counterparts (e.g. SVD, QR) and because the products can be computed in a straightforward way in parallel (see also [8]).

Non-negative tensor factorizations have been introduced in the literature recently as well, and just like the unconstrained counterparts, the type of decomposition used varies [2]. In [22, 4], the authors consider non-negative tensor decompositions based specifically on the $t$-product, the approaches in the two papers differing by the additional constraints on the optimization as well as the algorithms proposed to compute the factorization. In [22], the authors develop an Alternating Direction Method of Multipliers (ADMM) [1] method for producing a non-negative patch tensor dictionary from a single, high resolution training image with the end goal of using the dictionary in the context of X-ray CT image reconstruction. The authors showed the method was good at producing reconstructions even for missing data situations, and that it gave improvements over matrix-based patch dictionary learning since the reconstructions were sparser and less sensitive to regularization parameters.

In this paper, we consider two classical problems – (lossy) image compression and image deblurring. In both applications, the first stage is to learn a tensor patch dictionary from multiple images of the same class using the approach in [22], and hence we review that problem briefly. Our first new contribution deals with finding a non-negative representation under the tensor $t$-product [9] of any image given a tensor dictionary. We give theoretical results and concrete illustrations that demonstrate the superiority of a tensor-patch dictionary over the corresponding matrix case. Then, we show how the modified residual norm steepest descent (MRNSD), [15, 6] can be utilized for non-negative tensor coefficient recovery under the $t$-product. Additionally, we introduce sparsity-inducing regularization to the algorithm that can lead to compressed representations for images. Furthermore, we show that constraining our image to the non-negative patch dictionary representation can lead to a new effective debluring approach that is robust to certain model mismatches.

This paper is organized as follows. Section 2 is devoted to the introduction of background and notation. In Section 3, we describe the process of patchification of images to make a tensor representation, and review the dictionary learning approach presented in [22] which we will use to generate our tensor dictionaries. In Section 4, we investigate the power of the tensor-tensor product based representation of images vs. the traditional matrix-based approach. The choice of parameters such as patch and dictionary sizes are relative to quality, storage and computation time are also considered here. Following that discussion is the MRNSD algorithm for tensors in Section 5. Here, we also discuss the incorporation of coefficient sparsity constraints into MRNSD to allow for the compressed representation of images. In Section 6 we give a short introduction to the image deblurring problem, explain how to represent the unknown image in terms of the tensor dictionary, and discuss the restoration problem that needs to be solved for the tensor coefficients. Numerical results are contained in Section Section 7 and a discussion and list of future work is given in Section 8. Detailed derivations for some of the claims are left to the appendicies.

**2. Notation and preliminaries.** A *tensor* is a multidimensional array of data; a first-order tensor is a vector and a second-order tensor is a matrix. This paper focuses on third-order tensors (i.e., three-dimensional data), though much of the theory can be extended to higher-order tensors. We denote tensors with script letters.

Suppose $\mathcal{A}$ is an $\ell \times m \times n$ tensor. As depicted in Figure 1, we can divide the tensor in several directions. *Frontal slices*, denoted $\mathbf{A}^{(k)}$ for $k = 1, \ldots, n$, are $\ell \times m$ matrices which fix the third-dimension of $\mathcal{A}$. *Lateral slices*, denoted $\vec{\mathcal{A}}_j$ for $j = 1, \ldots, m$, are $\ell \times 1 \times n$ tensors which fix the second-dimension of $\mathcal{A}$; we consider lateral slices to be $\ell \times n$ matrices oriented along the third-dimension. *Tube fibers or tubes*, denote $\mathbf{a}_{ij}$ for $i = 1, \ldots, \ell$ and $j = 1, \ldots, m$, are the $1 \times 1 \times n$ mode-3 fibers of $\mathcal{A}$ or $n \times 1$ column vectors oriented along the third dimension.
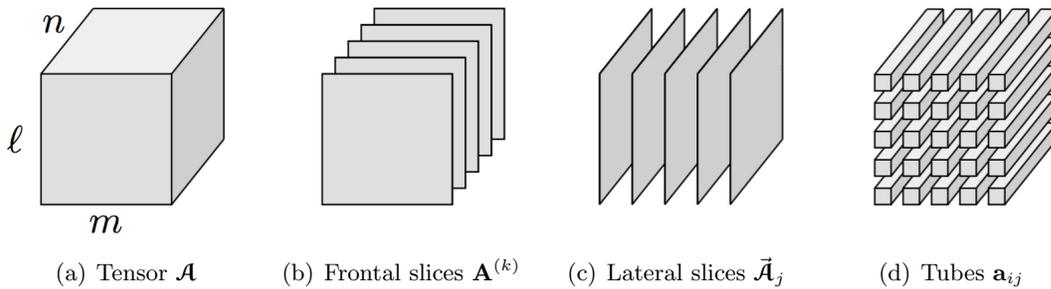


(a) Tensor $\mathcal{A}$  (b) Frontal slices $\mathbf{A}^{(k)}$  (c) Lateral slices $\vec{\mathcal{A}}_j$  (d) Tubes $\mathbf{a}_{ij}$

**Figure 1.** *Tensor notation.*

If we consider tensors as linear operators analogous to matrices, lateral slices are the analogous to column vectors, hence the notation $\vec{\mathcal{A}}$. In particular, tensors act on lateral slices just as matrices act on column vectors. Furthermore, lateral slices form the range and null space of a tensor [8]. For more detailed analysis on tensor linear algebra, we reference [8].

Many of the following definitions are taken directly from [9]. Using the $\ell \times m \times n$ tensor $\mathcal{A}$ illustrated in Figure 1, we define the `unfold` and `fold` operations as follows:

$$
(1) \qquad \mathtt{unfold}(\mathcal{A}) = \underbrace{\begin{pmatrix} \mathbf{A}^{(1)} \\ \mathbf{A}^{(2)} \\ \vdots \\ \mathbf{A}^{(n)} \end{pmatrix}}_{\ell n \times m}, \qquad \mathtt{fold}(\mathtt{unfold}(\mathcal{A})) = \mathcal{A}.
$$

The `unfold` function reshapes a tensor $\mathcal{A}$ into a block-column vector where each block is a frontal slice. The `fold` function reshapes an unfolded tensor into its original structure. Notice that the number of elements of $\mathcal{A}$ and $\mathtt{unfold}(\mathcal{A})$ is the same.

We now define the function `circ` which transforms a tensor $\mathcal{A}$ into a block-circulant matrix

whose blocks are the frontal slices of $\mathcal{A}$.

$$
(2) \qquad \texttt{circ}(\mathcal{A}) = \begin{pmatrix} \mathbf{A}^{(1)} & \mathbf{A}^{(n)} & \dots & \mathbf{A}^{(2)} \\ \mathbf{A}^{(2)} & \mathbf{A}^{(1)} & \dots & \mathbf{A}^{(3)} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{A}^{(n)} & \mathbf{A}^{(n-1)} & \dots & \mathbf{A}^{(1)} \end{pmatrix}.
$$

$$
\ell n \times mn
$$

Notice that the first column of $\texttt{circ}(\mathcal{A})$ is the unfolded tensor from Equation (1). Furthermore, notice that $\texttt{circ}(\mathcal{A})$ has $n$ times the number of elements of the original tensor $\mathcal{A}$. Fortunately, we need not form $\texttt{circ}(\mathcal{A})$ explicitly.

Using (1) and (2), the $t$-product of two tensors is defined in [9] as follows:

**Definition 2.1** ($t$-product). *Given $\mathcal{A}$ is an $\ell \times p \times n$ tensor and $\mathcal{B}$ is $p \times m \times n$, we define the $t$-product as*

$$
\mathcal{A} * \mathcal{B} = \texttt{fold}\left(\texttt{circ}(\mathcal{A}) \cdot \texttt{unfold}(\mathcal{B})\right),
$$

*where $\mathcal{A} * \mathcal{B}$ is an $\ell \times m \times n$ tensor and " $*$ " denote the $t$-product.*

Note that tubes commute under the $t$-product, and thus act analogously to scalars.

For the algorithms we describe in Section 5, we require the following two tensor norms [10, 22].

**Definition 2.2** (Frobenius norm). *Suppose $\mathcal{A}$ is an $\ell \times m \times n$ tensor. Then:*

$$
\|\mathcal{A}\|_F^2 = \texttt{trace}\left((\mathcal{A}^T * \mathcal{A})^{(1)}\right) = \sum_{k=1}^{n}\sum_{j=1}^{m}\sum_{i=1}^{\ell}(\mathcal{A}_{ij}^{(k)})^2.
$$

**Definition 2.3** (Sum norm). *Suppose $\mathcal{A}$ is an $\ell \times m \times n$ tensor. The sum norm is*

$$
\|\mathcal{A}\|_{\text{sum}} = \sum_{k=1}^{n}\sum_{j=1}^{m}\sum_{i=1}^{\ell}|\mathcal{A}_{ij}^{(k)}|.
$$

**2.1. Properties of the $t$-product.** As we alluded to above and given in [9], we can compute the $t$-product (see Definition 2.1) more efficiently using the Fourier transform:

**Definition 2.4** ($t$-product with Fourier transform). *Given $\mathcal{A}$ is an $\ell \times p \times n$ tensor and $\mathcal{B}$ is $p \times m \times n$, the $t$-product $\mathcal{C} = \mathcal{A} * \mathcal{B}$ can be computed as follows:*

$$
\widehat{\mathbf{C}}^{(i)} = \widehat{\mathbf{A}}^{(i)} \cdot \widehat{\mathbf{B}}^{(i)} \quad \text{for } i = 1, \dots n,
$$

*where $\widehat{\mathcal{A}} = \texttt{fft}(\mathcal{A}, [\,], 3)$, $\mathcal{C} = \texttt{ifft}(\widehat{\mathcal{C}}, [\,], 3)$, and $\texttt{fft}$, $\texttt{ifft}$ are the one-dimensional fast Fourier and inverse Fourier transforms, respectively, applied along the third-dimension.*

Definition 2.4 can be implemented in parallel perfectly, hence is an efficient algorithm for computing the $t$-product.

An alternative perspective on the $t$-product will be essential to our understanding of tensor dictionary learning in Subsection 4.2. (See also [4, 22].) Suppose $\vec{\mathcal{A}}$ is an $\ell \times 1 \times n$ lateral slice. Then, $\texttt{squeeze}(\vec{\mathcal{A}})$ rotates the lateral slice into an $\ell \times m$ matrix; the $\texttt{twist}$ transformation reverses this process (see Figure 2).
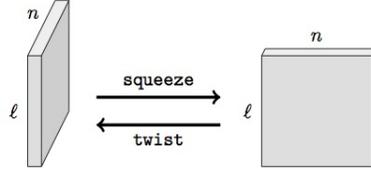


**Figure 2.** *Illustration of* $\texttt{squeeze}$ *and* $\texttt{twist}$ *transformations.*

Next we show how the structure imposed by the $t$-product impacts lateral slices.

Definition 2.5 ($t$-product with lateral slices). *Given $\mathcal{A}$ is an $\ell \times p \times n$ tensor and $\mathcal{B}$ is $p \times m \times n$, we can write the $\mathcal{C} = \mathcal{A} * \mathcal{B}$ as follows:*

$$\texttt{squeeze}\left(\vec{\mathcal{C}}_j\right) = \sum_{i=1}^{p} \texttt{squeeze}\left(\vec{\mathcal{A}}_i\right) \cdot \texttt{circ}(\mathbf{b}_{ij}^T) \quad \text{for } j = 1, \ldots, m.$$

Originally, we viewed the $t$-product as $\mathcal{A}$ acting on the lateral slices of $\mathcal{B}$ (see Definition 2.1). The significance of Definition 2.5, originally noted in [4], is that we can consider tubes of $\mathcal{B}$ to be "coefficients" of lateral slices of $\mathcal{A}$. We say more about this in Section 4.2.

**3. Patch Tensor Representation and Learning.** We briefly discuss the general idea of dictionary learning with tensors. For more background on matrix-based dictionary learning, one can see [21] and the references therein. As this paper focuses on tensor formulations, we keep our overview of the literature to describing the tools from [22] that we use here.

**3.1. Image to Tensor Mapping.** First, let us describe the transformation of a single two dimensional image into a third-order tensor. Let us suppose we have one image $\mathbf{B}$ of size $N_r \times N_c$, and we desire to consider this image in terms of $p \times q$ patches, where $N_r = pn_r$ and $N_c = qn_c$ for some integers $n_r, n_c$, respectively. Then our $N_r \times N_c$ image can be mapped to a $p \times n_r n_c \times q$ third order tensor $\mathcal{B}$ by putting each image patch into a lateral slice of our tensor. We choose to use a lexicographical ordering by patch columns. As seen in Figure 3, the (1,1) patch in $\mathbf{B}$ is mapped to the first lateral slice in $\mathcal{B}$ (i.e. $\mathcal{B}_{:,1,:}$), the (2,1) patch in $B$ becomes the 2nd lateral slice of $\mathcal{B}$, etc. Clearly, the process is completely reversible: given the patch tensor representation of an image, we can map back to its matrix representation.

In sum, $\mathbf{b}, \mathbf{B}, \mathcal{B}$ all represent the same image, but in different formats. Relative dimensions are summarized in Table 1 for easy reference.
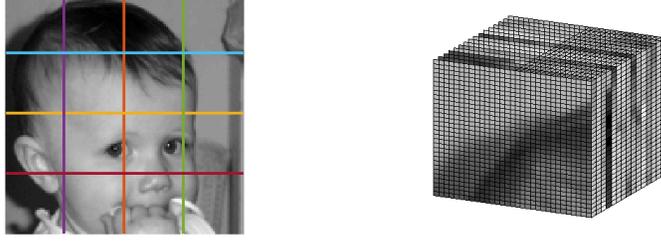
**Figure 3.** *Illustration of tensor patchification of image (left) to construct its tensor representation (right).*

| $N_r = n_r p$; $N_c = n_c q$; $M = n_r n_c$ | | | | |
|:---:|:---:|:---:|:---:|:---:|
| **B** | $\mathcal{B}$ | **b** | $\mathcal{D}$ | $\mathcal{C}$ |
| $N_r \times N_c$ | $p \times M \times q$ | $N_r N_c \times 1$ | $p \times s \times q$ | $s \times M \times q$ |

**Table 1**

*Left three columns give the dimensions of the various representations of the same image. In the right two columns, the image approximation in tensor form is assumed $\mathcal{X} = \mathcal{D} * \mathcal{C}$, and the corresponding sizes of $\mathcal{D}$ and $\mathcal{C}$ under this assumption are given.*

**3.2. Tensor-based dictionary learning.** Now suppose we have a sample space of $N_I$ images, each image of size $N_r \times N_c$. Following [22], we divide each image into patches of size $p \times q$ and let $M$ be the number of patches per image. Unlike matrix-based dictionary learning, we do not vectorize each patch. Instead, we store all patches as lateral slices of a sample space tensor $\mathcal{Y}$ of size $p \times t \times q$ where $t = N_I \cdot M$ is the total number of patches (see Figure 4).
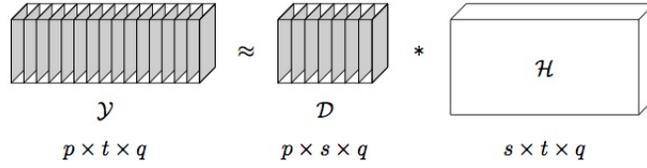


**Figure 4.** *Illustration of tensor dictionary learning decomposition.*

To 'learn' the dictionary representation is to minimize $\|\mathcal{Y} - \mathcal{D} * \mathcal{H}\|_F^2$ where $\mathcal{D} \in \mathbb{R}_+^{p \times s \times q}$, $\mathcal{H} \in \mathbb{R}_+^{s \times t \times q}$, and $s \ll t$. Here, $\mathcal{H}$ contains the tensor coefficients for the tensor dictionary $\mathcal{D}$. From [22], the problem to solve is

$$(3) \qquad \min_{\mathcal{D}, \mathcal{H}, \mathcal{U}, \mathcal{V}} \quad \frac{1}{2}\|\mathcal{Y} - \mathcal{U} * \mathcal{V}\|_F^2 + \lambda\|\mathcal{H}\|_{\text{sum}} + I_{\mathbb{R}_+^{s \times t \times q}}(\mathcal{H}) + I_{\mathsf{D}}(\mathcal{D})$$

$$\text{subject to} \quad \mathcal{D} = \mathcal{U} \quad \text{and} \quad \mathcal{H} = \mathcal{V},$$

where $\mathcal{D}, \mathcal{U} \in \mathbb{R}_+^{p \times s \times q}$ and $\mathcal{H}, \mathcal{V} \in \mathbb{R}_+^{s \times t \times q}$. In Equation (3), $\lambda$ is a regularization parameter and the sum norm (see Definition 2.3) promotes sparsity of the coefficient tensor $\mathcal{H}$. We denote the indicator function of a set $Z$ as $I_Z$. Thus, $I_{\mathbb{R}_+^{s \times t \times q}}(\mathcal{H})$ ensures the coefficients $\mathcal{H}$

are non-negative and $I_{\mathtt{D}}$ ensures the dictionary $\boldsymbol{\mathcal{D}}$ belongs to the compact and convex set (4)

$$(4) \qquad \mathtt{D} \equiv \left\{ \boldsymbol{\mathcal{D}} \in \mathbb{R}_+^{p \times s \times q} \mid \|\vec{\boldsymbol{\mathcal{D}}}_i\|_F \leq \sqrt{pq},\ i = 1, \ldots, s \right\}.$$

As described in [22], we impose the extra constraint that $\boldsymbol{\mathcal{D}} \in \mathtt{D}$ (as opposed to $\boldsymbol{\mathcal{D}} \in \mathbb{R}_+^{p \times s \times q}$) to avoid scaling ambiguity; that is, for any $\beta > 0$, $\|\boldsymbol{\mathcal{Y}} - (\beta \cdot \boldsymbol{\mathcal{D}}) * (\frac{1}{\beta}\boldsymbol{\mathcal{H}})\|_F^2 = \|\boldsymbol{\mathcal{Y}} - \boldsymbol{\mathcal{D}} * \boldsymbol{\mathcal{H}}\|_F^2$.

We will not discuss the specifics of the tensor-based ADMM algorithm in this paper; we refer the reader to [22] for a full analysis. When optimizing (3), we project $\boldsymbol{\mathcal{D}}$ and $\boldsymbol{\mathcal{H}}$ into $\mathtt{D}$ and $\mathbb{R}_+^{s \times t \times q}$, respectively. We choose to project $\boldsymbol{\mathcal{D}}$ into $\mathtt{D}$ using the infinity norm, that is:

$$(5) \qquad P_{\mathtt{D}}(\boldsymbol{\mathcal{D}})_{ij}^{(k)} = \min(\max(\boldsymbol{\mathcal{D}}_{ij}^{(k)}, 0), 1),$$

where $P_{\mathtt{D}}$ is the projection operator, an option included in the publically available code [20].

**3.3. Representation/Recovery Formulation.** Now suppose we have image $\mathbf{B} \in \mathbb{R}_+^{N_r \times N_c}$ which we would like to represent in terms of a dictionary we learn by solving (3). Let $\boldsymbol{\mathcal{B}}$ be the $p \times M \times q$ be the patchified tensor representation, and $\boldsymbol{\mathcal{D}}$ the $p \times s \times q$ non-negative patch dictionary. To represent the non-negative image via our patch dictionary, we solve

$$(6) \qquad \min_{\boldsymbol{\mathcal{C}}} \frac{1}{2}\|\boldsymbol{\mathcal{B}} - \boldsymbol{\mathcal{D}} * \boldsymbol{\mathcal{C}}\|_F^2 \quad \text{subject to } \boldsymbol{\mathcal{C}} \in \mathbb{R}_+^{s \times M \times q}.$$

In other words, if we can determine $\boldsymbol{\mathcal{C}} \geq 0$ such that $\boldsymbol{\mathcal{B}} \approx \boldsymbol{\mathcal{D}} * \boldsymbol{\mathcal{C}}$, then the image approximation is obtained by computing $\boldsymbol{\mathcal{D}} * \boldsymbol{\mathcal{C}}$ and mapping the resulting tensor back to a 2D image by inverting the patchification process. But several issues warrant discussion before presentation of the algorithm to solve for $\boldsymbol{\mathcal{C}}$. First, we need to give intuition as to why the tensor-based approach to the image model can provide significantly different results than the a matrix-based analogue, independent of the method produced to generate the dictionary. Then, we need to consider choices of patch and dictionary sizes required to maximize the representation power and harness the computational efficiencies of the tensor-based approach. These issues are covered in the next section.

**4. The Tensor Formulation: Advantages and Parameter Choices.** We first explain the power behind the tensor-based approach. Then, we discuss the choice of parameters such as dictionary size and patch size to maximize the potential of our new method.

**4.1. Tensor Superiority.** In this subsection, we will assume that a patch-dictionary has already been determined. It does not matter for the moment how that dictionary was derived: our goal is to show the differences in the solution sets to the two problems of image approximation, one based on a matrix-formulation, and one based on the tensor-formulation.

Let $\boldsymbol{\mathcal{D}} \in \mathbb{R}_+^{p \times s \times q}$ denote the dictionary in tensor form, and define $\underline{\mathbf{D}} = \mathtt{unfold}\,(\boldsymbol{\mathcal{D}}) \in \mathbb{R}_+^{pq \times s}$. Likewise, let $\boldsymbol{\mathcal{B}}$ denote the patchified tensor image, and let $\underline{\mathbf{B}} = \mathtt{unfold}\,(\boldsymbol{\mathcal{B}}) \in \mathbb{R}_+^{pq \times n_r n_c}$. We have the following theorem:

**Theorem 4.1.** *Consider the set of solutions to within a tolerance $\epsilon$:*

$$\mathcal{X}_{mat} := \{\mathbf{C} \in \mathbb{R}_+^{s \times n_r n_c} \| \|\underline{\mathbf{B}} - \underline{\mathbf{D}} \cdot \mathbf{C}\|_F \leq \epsilon\}$$

$$\mathcal{X}_{ten} := \{\mathcal{C} \in \mathbb{R}_+^{s \times n_r n_c \times q} | \|\mathcal{B} - \mathcal{D} * \mathcal{C}\|_F \leq \epsilon\}.$$

Let $\mathcal{X}_{mat,e}$ denote the set of tensors of size $s \times n_r n_c \times q$ whose first frontal slice is from $\mathcal{X}_{mat}$ and the remaining $q - 1$ frontal slices are zeros. Then $\mathcal{X}_{mat,e} \subset \mathcal{X}_{ten}$. That is, the set of solutions of the tensor problem effectively contains the set of solutions to the matrix problem.

*Proof.* See [16]. ∎

This suggests that in solving the tensor problem, the solutions to the matrix problem are achievable, and we would be able to recover those if those are optimal in the tensor framework, as we demonstrate in Example 4.1 below. However, the tensor case may provide better solutions by virtue of working in the tensor algebra, which we see in Example 4.2.

**Example 4.1.** *Let* $\mathbf{B} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$ *be a single patch ($p = q = 2$, $n_r = n_c = 1$) which we interpret as the entire image, meaning* $\underline{B} = \mathtt{vec}(\mathbf{B})$ *is* $4 \times 1$. *Suppose*

$$\underline{D} = \begin{bmatrix} 1 & 0 & 0 & 0 & \frac{1}{4} \\ 0 & 1 & 0 & 0 & \frac{1}{2} \\ 0 & 0 & 1 & 0 & \frac{3}{4} \\ 0 & 0 & 0 & 1 & 1 \end{bmatrix}.$$

*Set* $\epsilon = 0$; *that is, find exact solutions* $\mathbf{c} \in \mathbb{R}_+^{5 \times 1}$ *such that* $\|\underline{B} - \underline{D} \cdot \mathbf{c}\| = 0$. *It is easily checked that* $\mathbf{c}_a = [1, 1, 1, 1, 0]^T$ *and* $\mathbf{c}_b = [3/4, 1/2, 1/4, 0, 1]^T$ *with* $\|\mathbf{c}_a\|_1 = 4$ *and* $\|\mathbf{c}_b\|_1 = 5/2$ *are both exact solutions of the matrix optimization problem. It is also easy to see a non-negative solution cannot be obtained with fewer than four non-zero coefficients.*

*We can exactly capture these matrix solutions in the tensor framework. Let* $\mathcal{B} = \mathtt{twist}(\mathbf{B})$ *be the patch stored as a* $2 \times 1 \times 2$ *lateral slice and let* $\mathcal{D} = \mathtt{fold}(\underline{D})$ *be the equivalent tensor dictionary of size* $2 \times 5 \times 2$. *We are trying to find exact solutions* $\mathcal{C} \in \mathbb{R}_+^{5 \times 1 \times 2}$ *such that* $\|\mathcal{B} - \mathcal{D} * \mathcal{C}\| = 0$. *If we let* $\mathcal{C}_{:,:,1} = \mathbf{c}_a$ *and* $\mathcal{C}_{:,:,2} = \mathbf{0}$, *then it is easily seen that* $\mathcal{C}$ *is a solution of the tensor version of the problem – this solution is effectively the matrix solution, in tensor form. However, the coefficient tensor*

$$\mathcal{C}_{:,:,1} = [1/3, 0, 0, 0, 2/3]^T \text{ and } \mathcal{C}_{:,:,2} = [1/3, 0, 0, 0, 2/3]^T,$$

*is a solution to the tensor version of the problem with no matrix-based analogue, and we also observe* $\|\mathcal{C}\|_{\mathrm{sum}} = 2$. *Thus, the set of tensor solutions is bigger, and for the same number of non-zeros in the coefficients, we can get tensor solutions of smaller sum norm.*

**Example 4.2.** *Here, we let* $\mathbf{B} = \begin{bmatrix} 1 & 3 \\ 2 & 4 \end{bmatrix}$, *still assuming a single patch, and we assume*

$$\underline{D} = \begin{bmatrix} 1 & 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 1 & 0 & 1 \\ 0 & 1 & 0 & 0 & 1 \end{bmatrix},$$

*with* $\mathcal{D} = fold(\underline{D})$. *It is easy to verify that ther is no non-negative* $\mathbf{c}$ *that can exactly recover* $\mathtt{vec}(\mathbf{B})$ *in the matrix case. If we set* $\epsilon = 1/2$, *then one element of* $\mathcal{X}_{mat}$ *is* $\mathbf{c} =$

$[1/2, 7/2, 2, 0, 1/2]^T$ *with* $\|\mathbf{c}\|_1 = 13/2$. *However, with only four non-zeros entries in our tensor coefficients, we can resolve* $\boldsymbol{\mathcal{B}}$ *exactly: e.g.,* $\boldsymbol{\mathcal{C}}_{:,:,1} = [0, 1, 1, 0, 0]^T$ *and* $\boldsymbol{\mathcal{C}}_{:,:,2} = [0, 1, 3, 0, 0]^T$ *satisfies* $\boldsymbol{\mathcal{B}} = \boldsymbol{\mathcal{D}} * \boldsymbol{\mathcal{C}}$, *has four non-zeros, and has* $\|\boldsymbol{\mathcal{C}}\|_{sum} = 6$.

*This demonstrates that we can get a richer and possibly more accurate set of solutions for the tensor representation of the problem than for the matrix version.*

**4.2. Parameters.** In [22], the authors use a value for $s$ that is consistent with a matrix-based patch dictionary learning algorithm, and illustrate on some CT image examples that when keeping $s$ fixed, the tensor patch dictionary allows for better approximation, but otherwise, the choices of $s$, $p$ and $q$ are not further discussed. Here, we explain why $s \geq p$ is necessary to get good representations. We then discuss why $s \gg p$ is not advantageous from a storage perspective, and explain why $s = 2p$ is sufficient from a qualitative point of view.

**Dictionary Dimesion $s$.** In [22], $s$ was chosen to be a small multiple of the product $pq$. The reasoning for this was that in matrix patch dictionary learning, each patch is expressed as a vector and then approximated as a linear combination of the columns of the dictionary matrix. Since the dictionary matrix would have $pq$ rows, then the choice of $s \geq pq$ would be required to try to ensure a spanning set. However, taking $s$ this large for the tensor dictionary case is in fact not necessary for reasonably sized patches, as we now explain. Further, large $s$ is not a good choice in terms of computational efficiency, as we show later.

From Definition (2.5), each of the $N$ image patches is approximated as

$$(7) \qquad \mathbf{B}_j \approx \sum_{i=1}^{s} \mathbf{D}_i \mathtt{circ}(\mathbf{c}_{ij}), \qquad j = 1, \ldots, M,$$

where $\mathbf{B}_j = \mathtt{squeeze}\left(\vec{\boldsymbol{\mathcal{B}}}_j\right)$ and $\mathbf{D}_i := \mathtt{squeeze}\left(\vec{\boldsymbol{\mathcal{D}}}_i\right)$ are both in $\mathbb{R}_+^{p \times q}$ and $\mathbf{c}_{ij} = \boldsymbol{\mathcal{C}}_{i,j,:}^T$.

Postmultiplication of $p \times q$ matrix $\mathbf{D}_i$ by a $q \times q$ circulant generated by the tube $\mathbf{c}_{ij}$ can be written

$$\mathbf{D}_i \mathtt{circ}\left(\mathbf{c}_{ij}\right) = \mathbf{c}_{ij}^{(1)}\mathbf{D}_j + \mathbf{c}_{ij}^{(2)}\mathbf{D}_j\mathbf{Z} + \cdots + \mathbf{c}_{ij}^{(q)}\mathbf{D}\mathbf{Z}^{q-1},$$

where $\mathbf{Z}$ denotes the $q \times q$ circulant downshift matrix (i.e. $\mathbf{Z}^q = \mathbf{I}$).

Since each term in the sum (7) admits such an expansion, after regrouping we obtain

$$(8) \qquad \mathbf{B}_j \approx \sum_{i=1}^{s} \mathbf{c}_{ij}^{(1)}\mathbf{D}_i + \sum_{j=1}^{s} \mathbf{c}_{ij}^{(2)}\mathbf{D}_i\mathbf{Z} + \cdots + \sum_{j=1}^{s} \mathbf{c}_{ij}^{(q)}\mathbf{D}_i\mathbf{Z}^{q-1},$$

meaning the $p \times q$ non-negative patch $\mathbf{B}_j$ is described by a linear combination of $sq$, $p \times q$ non-negative matrices, although subsets of those matrices in the expansion are related via column shifts. We know that a spanning set for all $p \times q$ matrices would need to be of dimension $pq$. We do not know if the matrices in the above expression are all independent so we do not know if $p = s$ is sufficient, but certainly we do need $s \geq p$. We found in practice that it was sufficient to take $s$ a small multiple of $p$ as long as the patch sizes were not too large. Typically $s = 2p$ was all that was needed in our experiments to get reasonable representations. Though one might argue a larger value of $s$ may result in sparser coefficients, there is a trade-off with respect to the computational cost.

**Storage of Tensor Coefficients.** Storage of the original image requires storage of $N_r N_c$ pixel values. Storage of $\mathcal{D}$ (assuming it is dense, which it may not be) requires $pqs$ numbers, while storage of $\mathcal{C}$ requires $sq\frac{N_r}{p}\frac{N_c}{q} = s\frac{N_r}{p}N_c$ numbers. Thus, if $s = 2p$, storage of $\mathcal{C}$ *assuming that $\mathcal{C}$ is dense* requires $2N_rN_c$ numbers, twice the amount of storage of the image itself. For the deblurring application, we will not be concerned with this additional storage – the coefficients are a means to an end (namely, producing a high quality restoration). For the compression application, however, our goal will be to produce a $\mathcal{C}$ that is sparse, so that only non-zeros need to be stored.

**Patch Sizes.** There clearly must be a lower bound on the patch size: in the extreme with $p = q = 1$, $s = 1$, the dictionary is only one non-zero constant and we cannot have compression because then $\mathcal{C}$ is the image itself. Choosing patch sizes too small undermines the power of the representation in (8), and since the implementation of the algorithms utilizes FFT's of length $q$, there will be too much inefficiency if $q$ is very small (see further discussion in Subsection 5.1). If $p$ is too large, since we have shown we need $s \geq p$, we would have a high storage cost for the dictionary. After we discuss the algorithm, we will see that we are further constrained by the computational impact of the choice of patch sizes.

**4.3. Global Interpretation: Image Resolution vs. Patch Size.** To gain intuition, we observe that by placing all the image patches into position in the image, our tensor approximation is equivalent to the matrix representation

$$\mathbf{B} \approx \sum_{j=1}^{s}(\mathbf{I}_{n_r} \otimes \mathbf{D}_j)\begin{bmatrix} \texttt{circ}\left(\mathbf{c}_{j,1}\right) & \texttt{circ}\left(\mathbf{c}_{j,n_r+1}\right) & \cdots & \texttt{circ}\left(\mathbf{c}_{j,n_r(n_c-1)+1}\right) \\ \texttt{circ}\left(\mathbf{c}_{j,2}\right) & \texttt{circ}\left(\mathbf{c}_{j,n_r+2}\right) & \cdots & \texttt{circ}\left(\mathbf{c}_{j,n_r(n_c-1)+2}\right) \\ \vdots & \vdots & \vdots & \vdots \\ \texttt{circ}\left(\mathbf{c}_{j,n_r}\right) & \cdots & \cdots & \texttt{circ}\left(\mathbf{c}_{j,n_rn_c}\right) \end{bmatrix},$$

where each circulant block in the block matrix is of size $q \times q$, and there are $n_r$ block rows and $n_c$ block columns. Thus, the image has a expansion in terms of a structured global dictionary $(\mathbf{I} \otimes \mathbf{D}_j), j = 1, \ldots s$, although such an expansion is never computed explicitly. From this we see that the same dictionary can be used to reconstruct the same image at different resolutions. We illustrate this in the numerical results.

**5. MRNSD for tensors.** Since

$$\|\mathcal{B} - \mathcal{D} * \mathcal{C}\|_F^2 = \|\texttt{unfold}(\mathcal{B}) - \texttt{circ}(\mathcal{D})\texttt{unfold}(\mathcal{C})\|_F^2.$$

let us (implicitly) define

$$\mathbf{v} = \texttt{vec}(\texttt{unfold}(\mathcal{B})), \qquad \mathbf{c} = \texttt{vec}(\texttt{unfold}(\mathcal{C})), \qquad \text{and } \mathbf{D} = \mathbf{I} \otimes \texttt{circ}(\mathcal{D}).$$

The MRNSD algorithm ([15, 6]) was developed to solve $\min_{\mathbf{c}\geq 0}\|\mathbf{v} - \mathbf{D}\mathbf{c}\|_2$, so it can clearly be applied to our formulation. Of course it would be foolish to form $\mathbf{D}$ explicitly. In fact, we can use an equivalent and elegant formulation of each MRNSD step that uses all the tensor mechanics, and therefore only requires we have a routine that performs the $t$-product. The algorithm is given in Algorithm 1, and the details of the equivalence to this approach are given in Appendix A.

---

**Algorithm 1** MRNSD with $t$-product

---

1: **Input:** image $\mathcal{B}$, dictionary $\mathcal{D}$, initial estimate $\mathcal{C}_0$
2: Form gradient $\mathcal{G}_0 = -\mathcal{D}^T * (\mathcal{B} - \mathcal{D} * \mathcal{C}_0)$
3: **for** $k = 0, 1, 2, \ldots$ **do**
4:     $\mathcal{S}_k = \mathcal{C}_k \odot \mathcal{G}_k$                                                          {Form search direction (Appendix A)}
5:     $\theta_k = \texttt{trace}[(\mathcal{S}_k^T * \mathcal{G}_k)^{(1)}] / \| \underbrace{\mathcal{D} * \mathcal{S}_k}_{\mathcal{W}_k} \|_F^2$           {Determine optimal step size (Appendix A)}
6:     $\alpha_k = \max\{\theta_k, \min_{\mathcal{S}_{ij}^{(\ell)} > 0} (\mathcal{C}_k)_{ij}^{(\ell)} / (\mathcal{S}_k)_{ij}^{(\ell)}\}$           {Ensure step size preserves non-negativity}
7:     $\mathcal{C}_{k+1} = \mathcal{C}_k - \alpha_k \cdot \mathcal{S}_k$                                                          {Update coefficients}
8:     $\mathcal{G}_{k+1} = \mathcal{G}_k - \alpha_k \cdot \mathcal{D}^T * \mathcal{W}_k$                                                          {Update gradient}
9: **end for**

---

**5.1. Implementation Details.** Per iteration in Algorithm 1, the dominant costs are the two products $\mathcal{W}_k := \mathcal{D} * \mathcal{S}_k$ and $\mathcal{D}^T * \mathcal{W}_k$. Recall the $t$-product is computed by moving into the Fourier domain (i.e. computing $\widehat{\mathcal{D}}, \widehat{\mathcal{S}}_k$ and and their facewise matrix-matrix products). Some computations can be reused. We note that $\widehat{\mathcal{W}_k}$ need not be recomputed, since those entries are already known from computing $\mathcal{W}_k$ in the step size computation. Also, entries of $\widehat{\mathcal{D}^T}$ are known from entries of $\mathcal{D}$. So we need only to assess the costs of computing $\widehat{\mathcal{D}}, \widehat{\mathcal{S}}_k$, and the costs of doing the $q$, matrix-matrix products $\widehat{\mathcal{D}}^{(\ell)} \widehat{\mathcal{S}}_k^{(\ell)}$ and $\widehat{\mathcal{D}^T}^{(\ell)} \widehat{\mathcal{W}}^{(\ell)}$.

The computational cost for the FFTs is $O(s \cdot \log_2(q) \cdot (pq + \frac{N_r N_c}{p}))$, and the computational cost for the matvecs is $O(s N_r N_c)$.

Note that the cost of the matrix multiplications is independent of the patch size if we assume serial implementation. At the other extreme, for $q$ processors, each processor would compute a single matrix-matrix product at $s \frac{N_r N_c}{q}$ flops, so a larger value of $q$ is beneficial. Either way, $q$ should not be too small or the constant in front of the cost to perform length-$q$ FFTs will not be suitably amortized. We already observed $s \geq p$. If $s = kp$ for a small integer $k$, the total flop count in serial is $O(kp^2 \log_2(q)q + N_r N_c k(p + \log_2(q)))$. Note the cost grows more slowly for $q$, suggesting $p \leq q$ may be desirable. For sufficiently small fixed $k, p, q$, the cost grows as the number of unknowns in the image.

**5.2. Compression.** When we reconstruct images via tensor MRNSD, we tend to generate coefficients $\mathcal{C}$ which contain many small values. This is because of the efficient encoding of information inherent in the $t$-product discussed previously.

We introduce a sparsity regularization to our MRNSD minimization motivated by a proximal-operator framework [18]. We briefly outline the proximal gradient method; specific details can be found in [18]. Traditionally, proximal algorithms are a class of convex optimization techniques solving problems of the following form:

$$(9) \qquad\qquad \min_{\mathbf{c}} h(\mathbf{c}) \equiv f(\mathbf{c}) + g(\mathbf{c}),$$

where $f$ is smooth and convex and $g$ is simple and convex. For example, $f$ could be the $\ell_2$-norm (i.e., quadratic) and $g$ could be an $\ell_1$-regularization (i.e., piecewise-linear).

We cannot use a conventional gradient-descent algorithm in (9) because $g$ need not be

differentiable. Instead, we define the *proximal operator* of $g$ as follows:

$$(10) \qquad \text{prox}_g(\mathbf{y}) = \arg\min_{\mathbf{c}} \left\{ g(\mathbf{c}) + \frac{1}{2}\|\mathbf{c} - \mathbf{y}\|_2^2 \right\}.$$

The intuition behind (10) is to balance a point which minimizes a function $g$ that is close to another point $\mathbf{y}$. This interpretation gives rise to a two-step procedure to solve (9):

1. Minimize $f$ using gradient descent: $\mathbf{y} = \mathbf{c}_k - \alpha_k \nabla f(\mathbf{c}_k)$.
2. Find a nearby point which minimizes $g$: $\mathbf{c}_{k+1} = \text{prox}_g(\mathbf{y})$.

We can apply this proximal operator framework to MRNSD with regularization. For simplicity, we derive our method for matrix-vector products, with the understanding that we can translate this to tensor notation in our case, as we show at the end of this section.

Ideally, we use $\ell_1$-regularization to promote sparsity in our original problem:

$$(11) \qquad \min_{\mathbf{c} \geq 0} \frac{1}{2}\|\mathbf{b} - \mathbf{D}\mathbf{c}\|_F^2 + \lambda\|\mathbf{c}\|_1.$$

Using MRNSD, we incorporate the non-negativity constraint into the optimization using the mapping $\mathbf{c} = e^{\mathbf{z}}$. However, because $e^{\mathbf{z}}$ is strictly positive, simply regularizing $e^{\mathbf{z}}$ will not promote sparsity.

We incorporate the constraint into our function using the following mapping:

$$\mathbf{c} = e^{\mathbf{z}} - \epsilon\mathbf{1},$$

where $\mathbf{1}$ denotes the vector of all ones. This means $\mathbf{c}_i > -\epsilon$ and we will take $\epsilon \to 0$. We now minimize the unconstrained problem:

$$\min_{\mathbf{z}} = \underbrace{\tfrac{1}{2}\|\mathbf{b} - \mathbf{D}(e^{\mathbf{z}} - \epsilon)\|_F^2}_{f} + \underbrace{\lambda\|e^{\mathbf{z}} - \epsilon\|_1}_{g}$$

We compute the gradient of $f$ and the proximal operator of $g$ as follows:

$$\begin{aligned}
\nabla f &= e^{\mathbf{z}} \odot [-\mathbf{D}^T(\mathbf{b} - \mathbf{D}(e^{\mathbf{z}} - \epsilon))] \\
&= (\mathbf{c} + \epsilon) \odot [-\mathbf{D}^T(\mathbf{b} - \mathbf{D}\mathbf{c})], \quad \epsilon \to 0. \\
&= \mathbf{c} \odot [-\mathbf{D}^T(\mathbf{b} - \mathbf{D}\mathbf{c})]
\end{aligned}$$

This is the same gradient we had before. Next we consider the proximal operator:

$$\begin{aligned}
\text{prox}_g(\mathbf{y}) &= \arg\min_{\mathbf{c}} \left\{ \tfrac{1}{2}\|\mathbf{c} - \mathbf{y}\|_F^2 + \lambda\|\mathbf{c}\|_1 \right\} \\
&= \arg\min_{\mathbf{z}} \left\{ \tfrac{1}{2}\|e^{\mathbf{z}} - \epsilon - \mathbf{y}\|_F^2 + \lambda\|e^{\mathbf{z}} - \epsilon\|_1 \right\}
\end{aligned}$$

We solve this by computing the (sub)gradient and setting it equal to zero as follows:

$$\begin{aligned}
e^{\mathbf{z}} \odot (e^{\mathbf{z}} - \epsilon - \mathbf{y}) + \lambda e^{\mathbf{z}} \odot \text{sign}(e^{\mathbf{z}} - \epsilon) &= \mathbf{0}. \\
e^{\mathbf{z}} - \epsilon - \mathbf{y} + \lambda \cdot \text{sign}(e^{\mathbf{z}} - \epsilon) &= \mathbf{0}.
\end{aligned}$$

We can map this back to $\mathbf{c}$ as follows:

$$\mathbf{c} - \mathbf{y} + \lambda \cdot \mathtt{sign}(\mathbf{c}) = 0.$$

The solution to the above equation is exactly the soft-thresholding operator. Therefore, our MRNSD iteration with encorporated $\ell_1$-regularization is the following:

$$\mathbf{c}_{k+1} = G_{\alpha_k \cdot \lambda}[\mathbf{c}_k - \alpha_k \cdot \mathbf{c}_k \odot (-\mathbf{D}^T(\mathbf{b} - \mathbf{D}\mathbf{c}_k))],$$

where $G_\mu$ is the soft-thresholding operator:

$$G_\mu[c] = \begin{cases} c - \mu, & c > \mu \\ 0, & |c| < \mu \\ c + \mu, & c < -\mu. \end{cases}$$

Using the same observations as in the start of this section that allowed us to move from the matrix formulation to the tensor formulation, we arrive at the sparsity-constrained tensor-MRNSD formulation by changing Algorithm 1, Line 7 to the following:

(12) $$\mathcal{C}_{k+1} = G_{\alpha_k \cdot \lambda}[\mathcal{C}_k - \alpha_k \cdot \mathcal{S}_k].$$

We conclude this section by noting that we are not the first to consider augmentation of MRNSD iterates in order to encourage sparsity. In [2], the authors suggest applying a sparsity-type constraint to an MRNSD step. But the text was without mathematical justification, and we found in our examples that incorporation of their suggestion did little to promote sparsity.

**6. Deblurring.** We briefly review the standard model for image blurring/deblurring to set the stage for our tensor-based deblurring approach. For more background see [3].

The basic blurring model given assuming known image $\mathbf{x}_{true} = \mathtt{vec}(\mathbf{X}_{true}) \in \mathbb{R}_+^{N_r \times N_c}$, is

$$\mathbf{A}\mathbf{x}_{true} + \mathbf{n} = \mathbf{b},$$

where $\mathbf{A} \in \mathbb{R}^{N_r M_r \times N_c M_c}$ is a blurring operator whose singular values decay rapidly to 0, $\mathbf{n}$ is the unknown white noise vector and $\mathbf{b}$ is the blurred noisy image in vector form; that is, $\mathbf{b} = \mathtt{vec}(\mathbf{B})$ where $\mathbf{B}$ is $M_r \times M_c$. Since $\mathbf{A}$ and $\mathbf{b}$ are known but the noise is not, one might be tempted to ignore the noise, and compute the minimum-norm, least squares solution to $\mathbf{A}\mathbf{x} = \mathbf{b}$. However, the ill-conditioning of the operator renders the least squares solution worthless, since small singular values magnify the noise present in the data.

Algorithms for computing estimates $\mathbf{x} \approx \mathbf{x}_{true}$ in the presence of noise are called regularization methods. Iterative solvers, such as MRNSD, can be used as regularization methods. Consider applying MRNSD to

$$\min_{\mathbf{x} \geq 0} \|\mathbf{b} - \mathbf{A}\mathbf{x}\|_2.$$

It will produce sequences of iterates $\mathbf{x}_k$. Those iterates tend to exhibit semi-convergent behavior in that they will approximate the noise-free solution $\mathbf{x}_{true}$ with increasing $k$, up to a point. After a particular iteration, the method starts to fit the noise in $\mathbf{b}$ to solve the optimization problem, and the solution begins to resemble the noise-contaminated solution. If the

stopping parameter is picked before the contamination happens, the method is considered to be a regularization method.

In our approach, we require non-negativity of the image estimate in addition to the fact that the image be comprised from a learned, non-negative patch dictionary. In other words, we want our image estimate (expressed as a tensor, $\mathcal{X}$) to be given by $\mathcal{X} \approx \mathcal{D} * \mathcal{C}_k$ for $\mathcal{C}_k \geq 0$. MRNSD can be used to treat this problem. But first we need to show that it is possible to express the term on the right in the norm via matrix-vector products (but yet still employ the tensor format for computational efficiency during actual implementation). Consider the relationship between the two formats of the same image: $\mathbf{x}$ and $\mathcal{X}$. We see that

$$(13) \qquad\qquad \mathbf{x} = \mathbf{P}\mathtt{vec}(\mathtt{unfold}(\mathcal{X})) \approx \mathbf{P}\mathtt{vec}(\mathtt{unfold}(\mathcal{D} * \mathcal{C})),$$

for a permutation matrix $\mathbf{P}$. Now $\mathtt{unfold}(\mathcal{D} * \mathcal{C}) = \underline{\mathbf{D}}\,\mathtt{unfold}(\mathcal{C})$, where $\underline{\mathbf{D}} = \mathtt{circ}(\mathcal{D})$. So $\mathtt{vec}(\mathtt{unfold}(\mathcal{D} * \mathcal{C})) = (\mathbf{I} \otimes \underline{\mathbf{D}})\mathtt{vec}(\mathtt{unfold}(\mathcal{C}))$. Thus, we can apply MRNSD to solve

$$\min_{\mathcal{C} \geq 0} \|\mathbf{b} - (\mathbf{A}\mathbf{P}(\mathbf{I} \otimes \underline{\mathbf{D}}))\,\mathtt{vec}(\mathtt{unfold}(\mathcal{C}))\|_2$$

though in practice, we construct neither $\mathbf{P}$ nor $\mathbf{I} \otimes \underline{\mathbf{D}}$ explicitly, since all the necessary computations can be done with permutation indicies and $t$-products with $\mathcal{D}$ and $\mathcal{D}^T$. The computational cost of one iteration is dominated by matrix-vector products with $\mathbf{A}, \mathbf{A}^T$ and products with $\mathcal{D}, \mathcal{D}^T$, which as we saw previously, for sufficiently small values of $p, q$ is a small multiple of the number of unknowns in the image (and can be efficiently parallelized).

**7. Numerical Experiments.** We illustrate the power of the tensor dictionaries to represent images, both qualitatively and quantitatively. In all examples, we represent square images and the dimensions of the image and patches are powers of two.

**7.1. Power of Tensor Representations.** As discussed in Section 4 and Theorem 4.1, the tensor representations can exactly capture the matrix representations and there are a greater number of possible tensor representations. To illustrate the advantages of using a tensor representation, we compare the representations with either a learned tensor dictionary $\mathcal{D} \in \mathbb{R}_+^{16 \times 32 \times 16}$ against representations from a learned matrix dictionary $\mathbf{D} \in \mathbb{R}_+^{256 \times 512}$. In both cases, we use patches of size $16 \times 16$, either stored as lateral slices of $\mathcal{D}$ or as columns of $\mathbf{D}$. The number of dictionary elements in each case is twice the size of the first dimension (i.e., the dictionaries are equally over-complete). Both dictionaries were formed solving the ADMM formulation from images of faces in the CalTech101 database [11].

We form our representations in Figure 5 using 200 MRNSD iterations (Algorithm 1) and we start with a random, normalized initial guess.

In Figure 5, we see the tensor representation in 5(b) is better than the matrix representation in 5(c), both numerically and qualitatively. Thus, not only are there more possible the tensor representations (see Theorem 4.1), the representation we form is better. This is somewhat surprising as the number of coefficients (i.e., the representation ability) for both the tensor and matrix cases is the the same. More specifically, the sizes are the following: the tensor coefficients $\mathcal{C} \in \mathbb{R}_+^{32 \times 1024 \times 16}$ and the matrix coefficients $\mathbf{C} \in \mathbb{R}^{512 \times 1024}$ where 1024 is the number of patches in our original image $\mathbf{B}$. A potential reason for this improved representation is that the patches stored in the tensor dictionary $\mathcal{D}$ maintain some spatial
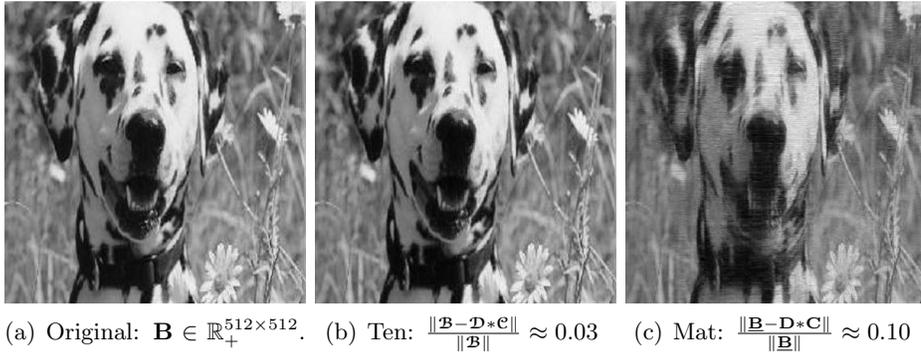
(a) Original: $\mathbf{B} \in \mathbb{R}_+^{512 \times 512}$.   (b) Ten: $\frac{\|\boldsymbol{\mathcal{B}} - \boldsymbol{\mathcal{D}} * \boldsymbol{\mathcal{C}}\|}{\|\boldsymbol{\mathcal{B}}\|} \approx 0.03$   (c) Mat: $\frac{\|\mathbf{B} - \mathbf{D} * \mathbf{C}\|}{\|\mathbf{B}\|} \approx 0.10$

**Figure 5.** *Comparison of tensor representation 5(b) with learned $\boldsymbol{\mathcal{D}} \in \mathbb{R}_+^{16 \times 32 \times 16}$ vs. learned matrix representation 5(c) with $\mathbf{D} \in \mathbb{R}_+^{256 \times 512}$.*

relationships typical in natural images (e.g., smooth curves) whereas the patches stored in the matrix dictionary $\mathbf{D}$ are more binary (e.g., sharp edges).

**7.2. Fixed Dictionary, Changing Resolution.** As noted previously, independent of how the dictionary was learned, we can employ that dictionary (assuming appropriate dimensions) on multiple resolutions of the same image, as illustrated in Figure 6. Importantly, the fact that the data was trained on images of a different resolution (in this case, the training data were all $128 \times 128$ images) is insignificant.
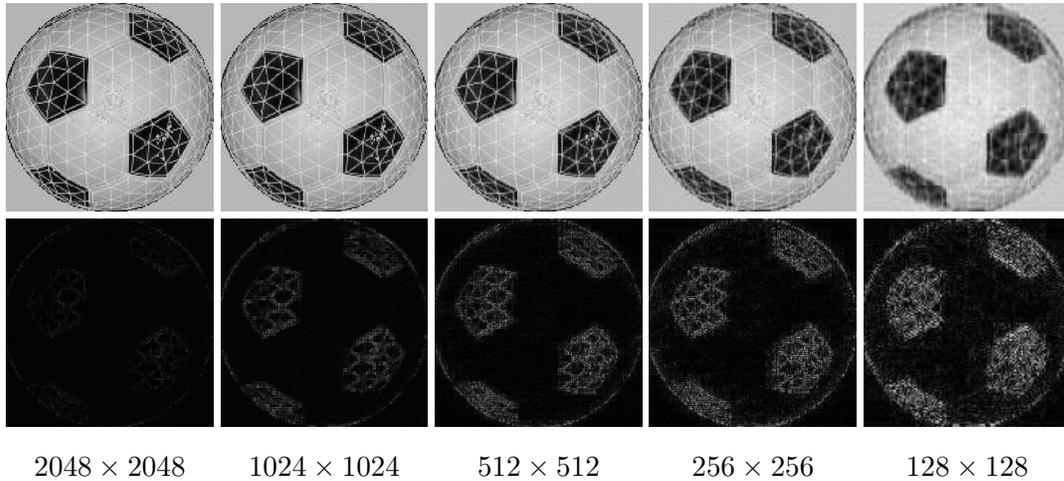


$2048 \times 2048$      $1024 \times 1024$      $512 \times 512$      $256 \times 256$      $128 \times 128$

**Figure 6.** *Effects of representations with $\boldsymbol{\mathcal{D}} \in \mathbb{R}_+^{16 \times 32 \times 16}$ as the image resolution changes. The dictionary was formed from $128 \times 128$ images. The top row shows the representations $\boldsymbol{\mathcal{D}} * \boldsymbol{\mathcal{C}}$ and the bottom row shows the absolute difference $|\boldsymbol{\mathcal{B}} - \boldsymbol{\mathcal{D}} * \boldsymbol{\mathcal{C}}|$. The resolution decreases from left to right.*

In Figure 6, we notice that as the image resolution decreases, the quality of our representations decreases as well (this is borne out by our reconstruction relative error). This is because the size of our patch relative to the image increases; i.e., each patch is representing a larger

portion of the image, and hence is less likely to match exactly. From another perspective, when we represent an image with a higher resolution, the dictionary patches act more like individual pixels in the image and hence provide a more accurate representation.

**7.3. Color.** We can also represent color images using the same dictionary generated from grayscale images. Suppose we have an RGB image of size $N_r \times N_c \times 3$ where the third dimension is the number of color channels. To patchify an RGB image, we treat each channel as a separate grayscale image from which we form patches and store as lateral slices of a tensor. This means we have three tensors of size $p \times M \times q$. We then concatenate the lateral slices of the patchified tensors for each color channel to obtain our RGB patchified tensor of size $p \times 3M \times q$. In Figure 7, we depict a representation of a color image using the same tensor dictionary $\boldsymbol{\mathcal{D}} \in \mathbb{R}_+^{16 \times 32 \times 16}$.



(a) Original $512 \times 512$, $\mathbf{B}$.    (b) Representation, $\boldsymbol{\mathcal{D}} * \boldsymbol{\mathcal{C}}$.    (c) Difference, $|\boldsymbol{\mathcal{B}} - \boldsymbol{\mathcal{D}} * \boldsymbol{\mathcal{C}}|$.

**Figure 7.** *Representing color images using* $\boldsymbol{\mathcal{D}} \in \mathbb{R}_+^{16 \times 32 \times 16}$. *The relative error of our representation is* $\|\boldsymbol{\mathcal{B}} - \boldsymbol{\mathcal{D}} * \boldsymbol{\mathcal{C}}\|/\|\boldsymbol{\mathcal{B}}\| \approx 0.02$.

**7.4. Compression.** In the examples, the $\mathbf{B}$ is $N \times N$. If we want to talk about the compression of a single image via the approximation $\boldsymbol{\mathcal{B}} \approx \boldsymbol{\mathcal{D}} * \boldsymbol{\mathcal{C}}$, we need to compute the compression ratio

$$\frac{\mathrm{nnz}(\boldsymbol{\mathcal{D}}) + \mathrm{nnz}(\boldsymbol{\mathcal{C}})}{N^2}.$$

However, if we are storing compressed representations of multiple images where they have all been compressed using the same dictionary, the cost of storing the dictionary becomes amortized over the multiple test images, so we approximate compression via $\mathrm{nnz}(\boldsymbol{\mathcal{C}})/N^2$.

For a fixed patch size, we know we want $s \geq p$. We can make $s$ larger (maybe a bit larger than $2p$), and increase sparsity to a point, but too big $s$ means too much non-uniqueness and the optimization problem gets trickier. We can change patch size. Increasing patch size for a fixed resolution beyond a certain point is not a good idea – we lose representability. But for larger images, we may well want to increase the patch size if we think our representation may be more sparse and we don't lose much representability. If we do that, $s$ must increase as a small multiple of $p$ and the cost of producing $\boldsymbol{\mathcal{C}}$ increases.

To examine the effects of patch size on compressibility, we compare the relative error to the approximate compression $\mathrm{nnz}(\boldsymbol{\mathcal{C}})/N^2$ where $N = 512$ in Figure 8. We use 200 MRNSD iterations (Algorithm 1) with the soft-thresholding step (12).
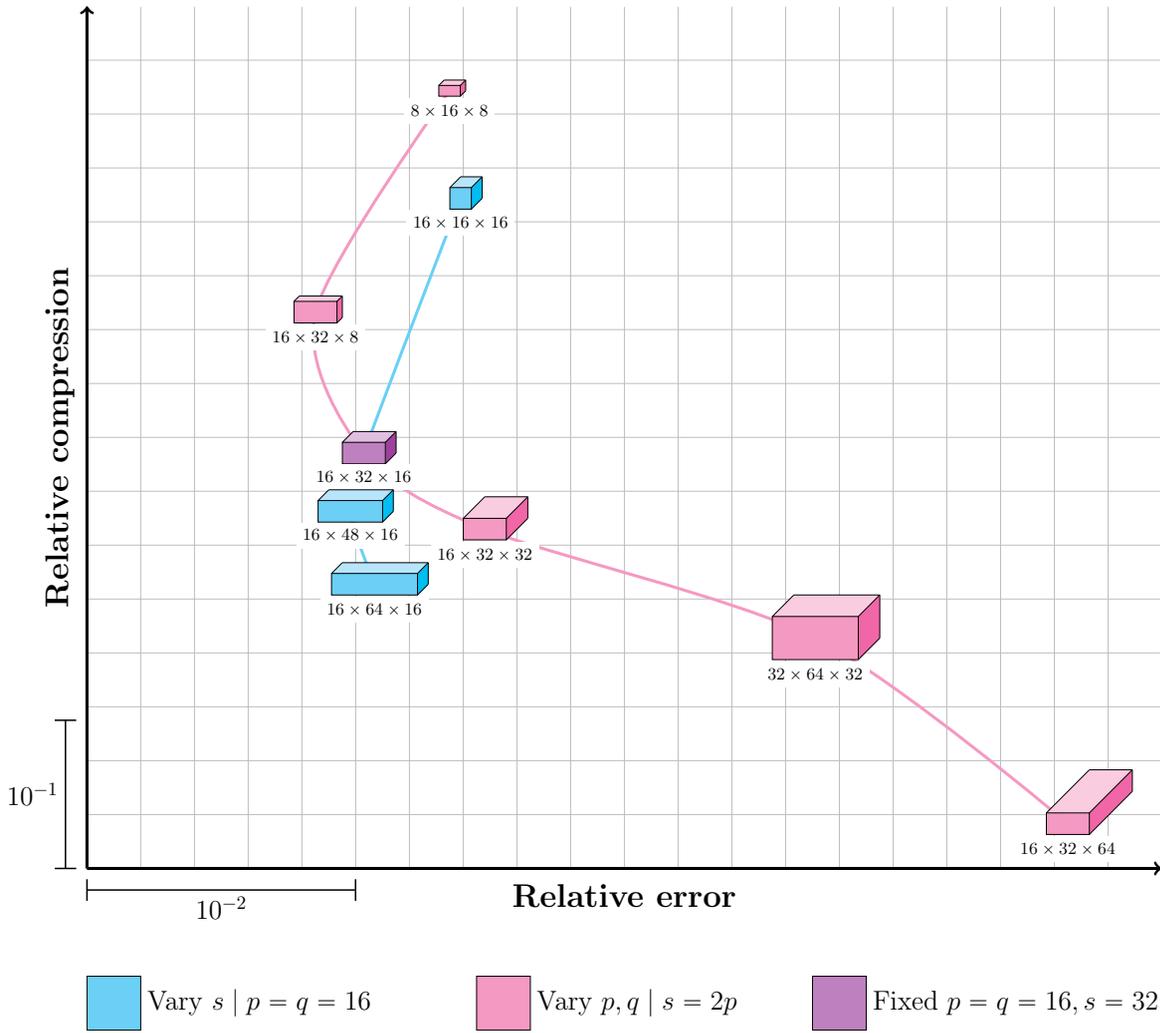
**Figure 8.** *Comparison of approximate compression for various dictionary sizes using sparsity-promoting MNRSD (Algorithm 1 with (12)) and a sparsity parameter of $\lambda = 10^{-10}$. The cyan line represents dictionaries of the same patch size, but varying the number of dictionary elements (i.e., width). The magenta line represents dictionary of various patch size, but the same level of of over-completeness (i.e., twice as many lateral slices as the first patch dimension p).*

There are a few key trends to notice in Figure 8. The first is that the more over-complete a dictionary is, the more compressed the representation without significant loss of accuracy (the cyan dictionaries). This behavior occurs because with a wider selection of dictionary patches to select, we likely need to select fewer patches to represent an image well. However, if we include the cost of storing these wider dictionaries, the compression ratio greatly increases due to the width of the dictionary.

The second trend is that the larger the patch size, the more compressed the representation, however with a significant loss of accuracy (the magenta dictionaries). This behavior occurs

because larger patches are able to capture larger sections of an image, hence fewer patches are required form the representation. However, the larger patches are less likely to reproduce the original image exactly, and hence decreases the representation quality. Interestingly, if we include the cost of storing the dictionaries with large patches, it does not significantly impact the overall storage cost – the width of the dictionary relative to the first patch size $p$ is a substantially more significant factor.

**7.5. Deblurring Results.** We used Matlab and features in the RestoreTools Matlab toolbox [14] as indicated.

**7.5.1. Example 1.** Our true image was the $512 \times 512$ images of the orca in Figure 9(a). We used the grain blur in Restoretools set of example files[1] to create a blurring operator $\mathbf{A}$ corresponding to reflexive boundary conditions. We computed $\mathbf{A}\mathbf{x}_{true}$, and added Gaussian noise at a noise level of 1 percent to the image. The blurred, noisy image in the figure.

Convergence to regularized solutions is known to be slow with MRNSD [15], so preconditioning is often used. Thus, in both the non-dictionary and dictionary reconstructions, we used the built-in preconditioner option and used MRNSD on the preconditioned problems

$$\min_{\mathbf{x} \geq 0} \|\mathbf{M}\mathbf{b} - \mathbf{M}\mathbf{A}\mathbf{x}\|_F \text{ or } \min_{\mathcal{C} \geq 0} \|\mathbf{M}\mathbf{b} - \mathbf{M}\left(\mathbf{A}\mathbf{P}(\mathbf{I} \otimes \underline{\mathbf{D}})\right)\texttt{vec}(\texttt{unfold}(\mathcal{C}))\|_F$$

where $\mathbf{M}$ denotes the preconditioner determined from the PSF and $\mathbf{b}$, using the default settings. The matrix $\mathbf{P}$ is a permutation matrix (see (13).

The algorithm needs a non-zero starting guess. In the matrix case, we used a vector of all ones as the initial guess for $\mathbf{x}$. In the tensor case, to make an equivalent comparison, we first formed a patchified version of an image of all ones. We then multiplied that by the tensor-pseudoinverse of $\mathfrak{D}$ (see [8] for details) and used this for the starting guess for $\mathcal{C}$.

We wanted to compare the quality of MRNSD with and without dictionaries. We do not discuss choosing optimal truncation parameters, though we note that the semi-convergence behavior is very much damped when using the dictionaries. We used two dictionaries derived from different data sets at different patch sizes. The first dictionary was obtained from the CalTech face database. We took $p = q = 16$ and $s = 32$. The second dictionary was obtained from a collection of 60 elephant photos [11]. Here, we took $p = q = 32$ and $s = 64$.

In the figure we compare the 'optimal' (i.e. solution at the iterate that gave smallest relative error against ground truth) solution with preconditioned MRNSD with no dictionary approach against other reconstructions. In Figure 9(c), we give the optimal reconstruction for the smaller dictionary. In Figure 9(d), we give the solution after 2000 iterations for the larger dictionary (the error is still decreasing at this point, so it may not be an optimal stopping point). In Figure 9(e), we averaged the solutions[2] to acknowledge the fact that this image has both fine scale features and components that are nearly uniform, so we expect that different resolution patches would be sensitive to this fact. We address the issue of multiresolution reconstructions in the Conclusions.

---

[1]The grain blur point-spread-function is for a 256 x 256 image, so we padded the blur by zeros to get a PSF suitable for a 512 x 512 image.

[2]In fact, any convex combination of the reconstructions would have been an option.

All dictionary based solutions gave reconstructions with smaller relative error and smaller structured similarity as shown in the table. It is worth noting that the dictionary-based reconstructions took longer to converge: while preconditioned MRNSD in the matrix-only case took 63 iterations to reach the optimal solution, it took the small dictionary 1,211 iterations, and as mentioned, we let the large dictionary case run 2000 iterations. On the other hand, this is not an entirely fair comparison, either, since the preconditioner was constructed relative to $\mathbf{A}$, whereas in the matrix-formulation of our tensor approach, we see the structure of the matrix-operator is quite a bit different.
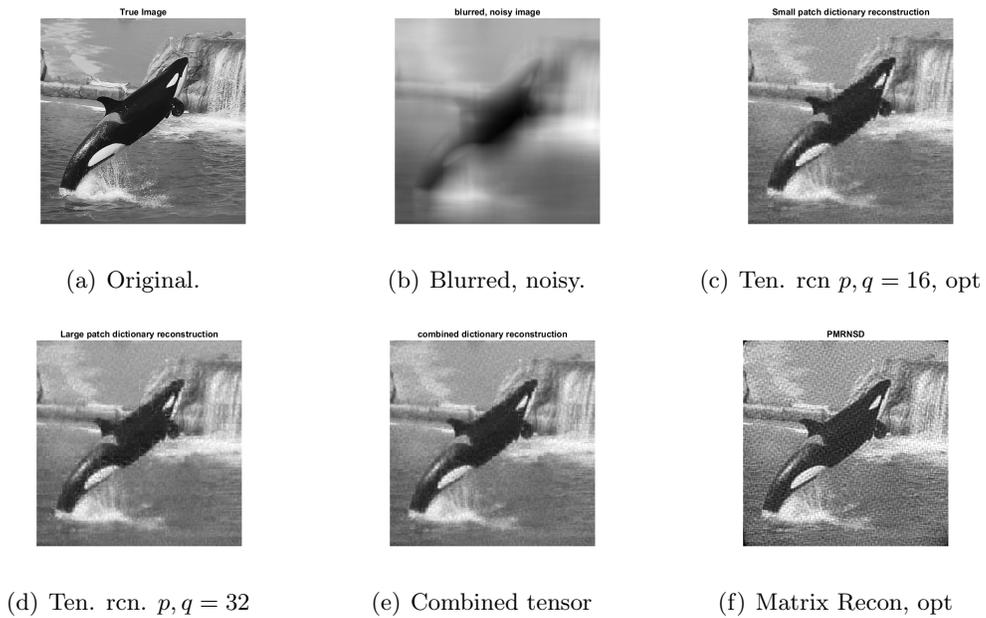


(a) Original.  (b) Blurred, noisy.  (c) Ten. rcn $p, q = 16$, opt

(d) Ten. rcn. $p, q = 32$  (e) Combined tensor  (f) Matrix Recon, opt

**Figure 9.** *Example 1: Orca original, blurred and noisy images, and various reconstruction results.*

|  | Small Dictionary | Large Dictionary | Combined | PMNRSD |
|---|---|---|---|---|
| Rel Err | 0.119 | 0.123 | 0.116 | 0.144 |
| SSIM | 0.518 | 0.522 | 0.541 | 0.376 |

**Table 2**

*Relative error and SSIM results for Example 1. The combined-dictionary image had slightly better relative error and SSIM results than any other. The matrix-based, non-dictionary reconstruction has notably worse relative error and SSIM to all the other dictionary-based reconstructions.*

**7.5.2. Examples 2 and 3: Underdetermined Problems.** In the first illustration, our true image was $256 \times 256$. We wanted to simulate a situation in which the boundary conditions of the blur were taken to be unknown. We took a symmetric Gaussian blur of discrete bandwidth 8 and $\sigma = 3$, and applied it to the true image, trimmed blurred the image by 8 pixels on all

sides, reshaped, and added 1 percent random Gaussian noise to the data[3] . This meant the data vector was only length $240^2$ while the true image was $256^2$, indicating there are fewer equations than unknowns.

Since the problem is underdetermined, it may be desirable to add regularization to enforce smooth transitions between patches. For an $N \times N$ image and $p \times p$ patches, we consider

$$\min_{\mathcal{C} \geq 0} \left\| \begin{bmatrix} \mathbf{b} \\ 0 \end{bmatrix} - \left( \begin{bmatrix} \mathbf{A} \\ \lambda \mathbf{L} \end{bmatrix} \mathbf{P}(\mathbf{I} \otimes \underline{\mathbf{D}}) \right) \text{vec}(\text{unfold}(\mathcal{C})) \right\|_2, \qquad \mathbf{L} = \begin{bmatrix} \mathbf{I} \otimes \mathbf{Q} \\ \mathbf{Q} \otimes \mathbf{I} \end{bmatrix},$$

where $\mathbf{Q}$ could either be an $(N-1) \times N$ first order discrete derivative operator, or, in order to minimize computation, an $(\frac{N}{p} - 1) \times N$ matrix approximating discrete derivatives only across patch jumps. We expect, for a suitable value of $\lambda$, some smoothing across patch boundaries. As our results below show, there is some modest gain that can be had when including extra regularization. However our method is relatively insensitive to choice of $\lambda$, whereas MRNSD without regularization is not.

The dictionary used in the reconstruction was constructed from the CalTech face data base (same as in the previous example). Patch sizes were $16 \times 16$ and we took $s = 32$ and used 2000 iterations to obtain each reconstruction (all convergence curves were nearly flat at this point). The original image in Figure 10(a) was obtained by cropping the MATLAB image `clutteredDesk.jpg`. The tensor dictionary based reconstructions for $\mathbf{L}$ a discrete gradient operator with $\lambda = 10$ is in Figure 10(c); the tensor-based reconstruction with no regularization is in Figure 10(d). The SSIM values of these were .815 and .776, respectively, showing the insensitivity to $\lambda$ and to additional regularization in general. The corresponding matrix MRNSD reconstructions with additional regularization (for $\lambda = 10$) and without additional regularization. Without regularization, the borders are white. The quality depends closely on the value of the regularization parameter, which is problematic. Moreover, the same or better quality reconstruction can be obtained using our tensor dictionary based approach without need of choosing a $\lambda$ – for example, the SSIM of the tensor-based reconstruction in (10(d)) was higher than for the matrix case for any value of $\lambda$ that we tried.

In the second illustration, the blurred and noisy image of size $512 \times 512$ is given in Figure 12. We use the same dictionary (i.e. learned from human faces) as in the previous illustration, but this time we used a discrete bandwidth of 12, $\sigma = 4$, and 5 percent Gaussian noise. In Figure 11, we show the relative errors for using our tensor approach with $\lambda = 100$ and patch-smoothing regularizer, the tensor approach with $\lambda = 0$ (i.e. no smoothing across patches), and the matrix-based MRNSD. We observe that the behavior is similar for the two tensor classes, with a slight improvement in the error observed when using the patch-regularization term. We note that semi-convergence behavior is observed in the matrix-based case whereas it is not observed in the tensor cases over the first 2000 iterations. We show the matrix-based reconstruction at the 'optimal' iteration count (198) in subfigure 12(e), and the reconstruction after 2000 iterations in subfigure 12(f). Even the optimal reconstruction is qualitatively not as good - clearly, there is a white boundary where it could not be reconstructed, and also

---

[3]To implement this process in MATLAB: Let $\mathbf{v} = \exp(-\frac{1}{\sqrt{2}\sigma}[0:7].^2)$, and $\mathbf{A}_1 = \text{toeplitz}(\mathbf{v})$. Define $\mathbf{T} = \mathbf{A}_1 * \mathbf{X}_{true} * \mathbf{A}_1$, and $\mathbf{b}_{true} = \mathbf{T}(8:247, 8:247)$, $\mathbf{b}_{true} = \mathbf{b}_{true}(:)$; and $\mathbf{b} = \mathbf{b}_{true} + c \cdot \text{randn}(\text{length}(\mathbf{b}_{true}, 1))$ where $c$ is such that the noise level is 0.01.

(a) Original.  (b) Blurred, noisy.  (c) Ten rcn $\lambda = 10$

(d) Ten rcn $\lambda = 0$  (e) Matrix rcn, $\lambda = 10$  (f) Matrix rcn, $\lambda = 0$
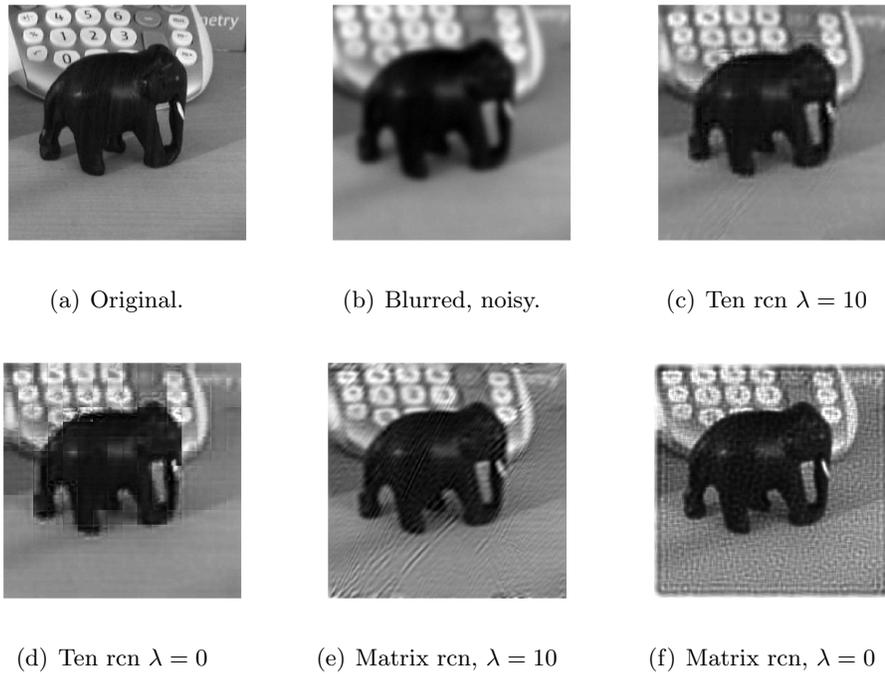
**Figure 10.** *Example 2: Original, blurred noisy image, and various reconstructions, with and without discrete derivative regularization for the tensor-dictionary-based and matrix reconstructions.*

there are ringing and fine scale noise artifacts in other areas of the matrix-based image as well. However, details are recovered using the tensor patch-based dictionary.

In Figure 12(d) we show a reconstruction using a different dictionary, also constructed from face data, of size $32 \times 64 \times 32$ for $\lambda = 100$ and 2000 iterations. We see that the quality is very close to the $16 \times 16$ patch dictionary, and there are improvements in some areas of the image but subtle degradataion in others. As we note in the conclusions, this suggests using a multilevel dictionary approach may improve the situation further.
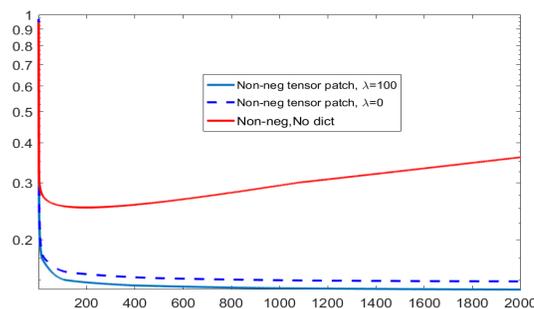


**Figure 11.** *Example 3: Convergence behavior. Horizontal axis is iteration number, vertical is relative error in the iterate against the true image.*

(a) Blurred, Noisy.       (b) Ten $\lambda = 100, p, q = 16$       (c) Ten $\lambda = 0, p, q = 16$

(d) Ten $\lambda = 100, p, q = 32$       (e) Opt. Matrix Recon       (f) Matrix Recon, 2000 its
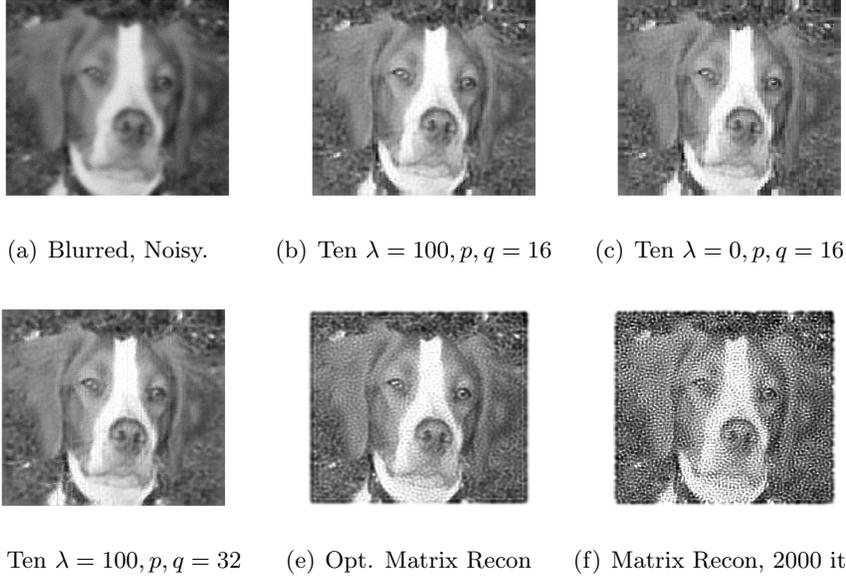
**Figure 12.** *Example 3. blurred and noisy image, several reconstructions. In the matrix cases, note the white border in the reconstructions due to the lack of information near the boundary.*

**8. Summary and Future Work.** In this work, we have shown the utility of learned tensor-patch dictionaries in the context of non-negative image representation, compression and image deblurring applications. In all cases, once a non-negative tensor patch-dictionary is available, we showed that the problems of compression and deblurring could be formulated in terms of recovering the corresponding non-negative tensor coefficient object. We gave an MRNSD tensor algorithm for finding the coefficient tensor, and described a modification that encourages sparsity in the coefficient tensor. Notably, this sparsity constraint is applicable whether or not one uses matrices or tensors in the formulation, thereby indicating the proposed approach has broader utility than for the purpose described here. In the case of deblurring, we showed the tensor representation is particularly effective in mitaging the effects of noise on the solution, especially in the case of underdetermined problems and boundary effects.

Importantly, we demonstrated that the class of data on which the dictionary is trained is surprisingly irrelevant in the context of image representation under the tensor-dictionary formulation, as is the resolution of the training data, both in the context of image compression and image deblurring. We also discussed issues related to patch size, and the trade-offs between sparsity, representability and computation time. We showed that a fixed dictionary can do remarkably well on representing images at various resolutions, and even across color channels. In our deblurring examples, we saw that the tensor dictionaries could mitigate semi-convergence behavior. We also observed that better representations could be obtained by convex combination of deblurred images constructed using dictionaries at different resolutions, which suggests further work is needed to design a multi-level dictionary representation that allows for better local feature description. Finally, as noted in [13, 22], the t-product generalizes to tensors of order higher than three, so our ideas generalize to higher order.

Recently, in [7], new tensor-tensor products have been defined which, like the t-product, permit a linear algebraic-type framework. The tensor-patch dictionary learning and representation approach can therefore be extended to these tensor-tensor products. Some of the preliminary details are offered in [16]. Further investigation into which class of tensor-tensor products provide for the best non-negative dictionaries for use in image compression and representation still needs to be considered, and is the subject of future research.

### Appendix A. Tensor-based MRNSD derivation.

Suppose $\mathcal{C} = e^{\mathcal{Z}}$ meaning $\mathcal{C}_{ij}^{(k)} = e^{\mathcal{Z}_{ij}^{(k)}}$. We then compute the search direction $\mathcal{S}$ by computing the gradient of 6 as follows:

$$(14) \qquad \begin{aligned} \mathcal{S} &= \nabla_{\mathcal{Z}} \left( \frac{1}{2} \|\mathcal{B} - \mathcal{D} * e^{\mathcal{Z}}\|_F^2 \right) \\ &= e^{\mathcal{Z}} \odot (-\mathcal{D}^T * (\mathcal{B} - \mathcal{D} * e^{\mathcal{Z}})). \end{aligned}$$

The search direction $\mathcal{S}$ is exactly the gradient of (6) with the addition of a Hadamard product $\odot$ with $\mathcal{C}$.

To determine the optimal step size, we solve for $\alpha$ as follows. For notational simplicity, we define the residual tensor $\mathcal{R}$, the gradient tensor $\mathcal{G}$, and $\mathcal{U}$ as follows:

$$\mathcal{R} = \mathcal{B} - \mathcal{D} * \mathcal{C}$$
$$\mathcal{G} = -\mathcal{D}^T * \mathcal{R}$$
$$\mathcal{U} = \mathcal{D} * \mathcal{S}$$

Note that $\mathcal{U}^T * \mathcal{R} = -\mathcal{S}^T * \mathcal{G}$. We reformulate (6) using Definition 2.2 in terms of $\mathcal{R}$, $\mathcal{G}$, and $\mathcal{U}$ as follows:

$$\begin{aligned} \frac{1}{2} \|\mathcal{B} - \mathcal{D} * (\mathcal{C} - \alpha \cdot \mathcal{S})\|_F^2 &= \frac{1}{2} \|\mathcal{R} + \alpha \cdot \mathcal{U}\|_F^2 \\ &= \frac{1}{2} \mathtt{trace}[((\mathcal{R} + \alpha \cdot \mathcal{U})^T * (\mathcal{R} + \alpha \cdot \mathcal{U}))^{(1)}] \\ &= \frac{1}{2} \mathtt{trace}[(\mathcal{R}^T * \mathcal{R} + 2\alpha \cdot \mathcal{U}^T * \mathcal{R} + \alpha^2 \cdot \mathcal{U}^T * \mathcal{U})^{(1)}]. \end{aligned}$$

Note that typically $\mathcal{U}^T * \mathcal{R} \neq \mathcal{R}^T * \mathcal{U}$; however, the trace of the first frontal slice is always equal. We made use of this fact in the last line above.

Now, we solve for $\alpha$ as follows:

$$\nabla_\alpha \frac{1}{2} \mathtt{trace}[(\mathcal{R}^T * \mathcal{R} + 2\alpha \cdot \mathcal{U}^T * \mathcal{R} + \alpha^2 \cdot \mathcal{U}^T * \mathcal{U})^{(1)}] = 0$$
$$\mathtt{trace}[(\mathcal{U}^T * \mathcal{R} + \alpha \cdot \mathcal{U}^T * \mathcal{U})^{(1)}] = 0$$

Solving for $\alpha$ and rewriting in terms of $\mathcal{D}$, $\mathcal{S}$, and $\mathcal{G}$, we get the optimal step size:

$$\begin{aligned} \alpha &= -\mathtt{trace}[(\mathcal{U}^T * \mathcal{R})^{(1)}]/\mathtt{trace}[(\mathcal{U}^T * \mathcal{U})^{(1)}] \\ &= \mathtt{trace}[(\mathcal{S}^T * \mathcal{G})^{(1)}]/\|\mathcal{D} * \mathcal{S}\|_F^2. \end{aligned}$$

We add an additional constraint on the $\alpha$ to ensure that we never move too far along the search direction and turn some coefficients $\mathcal{C}$ to negative values.

$$\theta = \mathtt{trace}[(\mathcal{S}^T * \mathcal{G})^{(1)}]/||\mathcal{D} * \mathcal{S}||_F^2$$
$$\alpha = \min\{\theta, \min_{\mathcal{S}_{ij}^{(k)}>0} (\mathcal{X}_{ij}^{(k)}/\mathcal{S}_{ij}^{(k)})\}.$$

## Appendix B. Quasiconvexity of MRNSD.

From Boyd and Vandenberghe's Convex Optimization, we have the following definition:

Definition B.1 (Quasiconvex).

*A function $f : \mathbb{R}^n \to \mathbb{R}$ is* quasiconvex *if all of its sublevel sets $S_\alpha = \{\mathbf{x} \in \mathbb{R}^n \mid f(\mathbf{x}) \leq \alpha\}$ for $\alpha \in \mathbb{R}$ are convex.*

We first expand $\Phi$ as follows:

$$\Phi(\mathbf{z}) = \tfrac{1}{2}\|\mathbf{D}e^{\mathbf{z}} - \mathbf{b}\|_F^2$$
$$= \tfrac{1}{2}\|\mathbf{D}e^{\mathbf{z}}\|^2 + \tfrac{1}{2}\|\mathbf{b}\|_F^2 - \mathbf{b}^T\mathbf{D}e^{\mathbf{z}}$$

Suppose for some $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ and $\alpha \in \mathbb{R}$, $\mathbf{s}, \mathbf{y} \in S_\alpha$; that is, $\Phi(\mathbf{x}), \Phi(\mathbf{y}) \leq \alpha$. We show that $\theta\mathbf{x} + (1-\theta)\mathbf{y} \in S_\alpha$ for all $\theta \in (0,1)$, hence that $S_\alpha$ is convex and $\Phi$ is quasiconvex.

$$\begin{aligned}
\Phi(\theta\mathbf{x} + (1-\theta)\mathbf{y}) &= \frac{1}{2}\|\mathbf{D}e^{\theta\mathbf{x}+(1-\theta)\mathbf{y}}\|_F^2 + \frac{1}{2}\|\mathbf{b}\|_F^2 - \mathbf{b}^T\mathbf{D}e^{\theta\mathbf{x}+(1-\theta)\mathbf{y}} \\
&\leq \frac{\theta}{2}\|\mathbf{D}e^{\mathbf{x}}\|_F^2 + \frac{1-\theta}{2}\|\mathbf{D}e^{\mathbf{y}}\|_F^2 + \frac{1}{2}\|\mathbf{b}\|_F^2 - \mathbf{b}^T\mathbf{D}e^{\theta\mathbf{x}+(1-\theta)\mathbf{y}} \quad \text{by convexity} \\
&\leq \theta\alpha + (1-\theta)\alpha - \mathbf{b}^T\mathbf{D}e^{\theta\mathbf{x}+(1-\theta)\mathbf{y}} \quad\quad\quad \text{by assumption} \\
&= \alpha - \mathbf{b}^T\mathbf{D}e^{\theta\mathbf{x}+(1-\theta)\mathbf{y}}
\end{aligned}$$

For our dictionary-learning problem, we assume $\mathbf{D}$ and $\mathbf{b}$ are non-negative because they are composed of images. Furthermore, $e^{\mathbf{z}}$ has non-negative components. Therefore,

$$\Phi(\theta\mathbf{x} + (1-\theta)\mathbf{y}) \leq \alpha - \mathbf{b}^T\mathbf{D}e^{\theta\mathbf{x}+(1-\theta)\mathbf{y}} \leq \alpha.$$

Therefore, $\theta\mathbf{x} + (1-\theta)\mathbf{y} \in S_\alpha$ and $S_\alpha$ is convex. Because $\Phi$ is quasiconvex, gradient descent will make progress towards a minimum (i.e., we will not be stuck at a saddle point).

### REFERENCES

[1] S. BOYD, N. PARIKH, E. CHU, B. PELEATO, AND J. ECKSTEIN, *Distributed optimization and statistical learning via the alternating direction method of multipliers*, Foundations and Trends in Machine Learning, 3 (2010), pp. 1–122.
[2] A. CICHOCKI, R. ZDUNEK, A. H. PHAN, AND S.-I. AMARI, *Nonnegative Matrix and Tensor Factorizations: Applications to Exploratory Multi-way Data Analysis and Blind Source Separation*, Wiley, 2009, https://pdfs.semanticscholar.org/94cc/6daad548a03c6edb0351d686c2d4aa364634.pdf.

[3] P. C. Hansen, J. G. Nagy, and D. P. O'Leary, *Deblurring Images: Matrices, Spectra, and Filtering (Fundamentals of Algorithms)*, SIAM, 2006.

[4] N. Hao, M. E. Kilmer, K. Braman, and R. C. Hoover, *Facial recognition using tensor-tensor decompositions*, SIAM Journal of Imaging Sciences, 6 (2013), pp. 457–463.

[5] B. Hunyadi, P. Dupont, W. Van Paesschen, and S. Van Huffel, *Tensor decompositions and data fusion in epileptic electroencephalography and functional magnetic resonance imaging data*, Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery, 7 (2017).

[6] L. Kaufman, *Maximum likelihood, least squares, and penalized least squares for pet*, IEEE Transactions on Medical Imaging, 12 (1993).

[7] E. Kernfeld, M. Kilmer, and S. Aeron, *Tensor–tensor products with invertible linear transforms*, Linear Algebra and its Applications, 485 (2015), pp. 545–570.

[8] M. E. Kilmer, K. Braman, N. Hao, and R. C. Hoover, *Third-order tensors as operators on matrices: a theoretical and computational framework with applications in imaging*, SIAM Journal on Matrix Analysis and Applications, 34 (2012), pp. 148–172.

[9] M. E. Kilmer and C. D. Martin, *Factorization strategies for third-order tensors*, Linear Algebra and its Applications, 435 (2011), pp. 641–658.

[10] T. Kolda and B. Bader, *Tensor decompositions and applications*, SIAM Review, 51 (2009), pp. 455–500.

[11] R. F. L. Fei-Fei and P. Perona, *Learning generative visual models from few training examples: an incremental bayesian approach tested on 101 object categories.*, IEEE CVPR Workshop on Generative-Model Based Vision., (2004).

[12] J. Liu, P. Musialski, P. Wonka, and J. Ye, *Tensor completion for estimating missing values in visual data*, Pattern Analysis and Machine Intelligence, IEEE Transactions on, 35 (2013), pp. 208–220.

[13] C. D. Martin, R. Shafer, and B. LaRue, *An order-p tensor factorization with applications in imaging*, SIAM Journal on Scientific Computing, 35 (2013), pp. A474–A490.

[14] J. Nagy, S. Berisha, J. Chung, K. Palmer, L. Perrone, and R. Wright, *RestoreTools: An object oriented matlab package for image restoration*, 2012, www.mathcs.emory.edu/~nagy/RestoreTools/index.html.

[15] J. Nagy and Z. Strakos, *Enforcing nonnegativity in image reconstruction algorithms*, Mathematical Modeling, Estimation, and Imaging, (2000).

[16] E. Newman, *A Step in the Right Dimension: Tensor Algebra and Applications*, PhD thesis, Tufts University, 2019.

[17] E. Newman, M. E. Kilmer, and L. Horesh, *Image classification using local tensor singular value decompositions*, in Proceedings from CAMSAP 2017, IEEE. See also arXiv preprint arXiv:1706.09693.

[18] N. Parikh and S. Boyd, *Proximal algorithms*, Foundations and Trends in Optimization, 1 (2013), pp. 123–231.

[19] O. Semerci, N. Hao, M. E. Kilmer, and E. L. Miller, *Tensor-based formulation and nuclear norm regularization for multienergy computed tomography*, IEEE Transactions on Image Processing, 23 (2014), pp. 1678–1693.

[20] S. Soltani, *Dlct-toolbox, a matlab package for the dictionary learning approach to tomograhic image reconstruction*, 2015, http://www.imm.dtu.dk/~pcha/HDtomo/.

[21] S. Soltani, M. Andersen, and P. C. Hansen, *Tomographic image reconstruction using training images*, Journal of Computational and Applied Mathematics, 313 (2016).

[22] S. Soltani, M. Kilmer, and P. C. Hansen, *A tensor-based dictionary learning approach to tomographic image reconstruction*, Spring BIT Numerical Mathematics, (2016), https://doi.org/10.1007/s10543-016-0607-z.

[23] M. A. O. Vasilescu and D. Terzopoulos, *Multilinear image analysis for facial recognition*, in Pattern Recognition, 2002. Proceedings. 16th International Conference on, vol. 2, IEEE, 2002, pp. 511–514.

[24] Z. Zhang, G. Ely, S. Aeron, N. Hao, and M. Kilmer, *Novel methods for multilinear data completion and de-noising based on tensor-svd*, in Proceedings of the 2014 IEEE Conference on Computer Vision and Pattern Recognition, CVPR '14, Washington, DC, USA, 2014, IEEE Computer Society, pp. 3842–3849, https://doi.org/10.1109/CVPR.2014.485, http://dx.doi.org/10.1109/CVPR.2014.485. Also excepted for oral presentation, 5.45 percent acceptance rate.