

# A Local Block Coordinate Descent Algorithm for the CSC Model

Ev Zisselman  
 Technion  
 Israel Institute of Technology  
 ev.zis@campus.technion.ac.il

Jeremias Sulam  
 Johns Hopkins University  
 jsulam1@jhu.edu

Michael Elad  
 Technion  
 Israel Institute of Technology  
 elad@cs.technion.ac.il

## Abstract

*The Convolutional Sparse Coding (CSC) model has recently gained considerable traction in the signal and image processing communities. By providing a global, yet tractable, model that operates on the whole image, the CSC was shown to overcome several limitations of the patch-based sparse model while achieving superior performance in various applications. Contemporary methods for pursuit and learning the CSC dictionary often rely on the Alternating Direction Method of Multipliers (ADMM) in the Fourier domain for the computational convenience of convolutions, while ignoring the local characterizations of the image. In this work we propose a new and simple approach that adopts a localized strategy, based on the Block Coordinate Descent algorithm. The proposed method, termed Local Block Coordinate Descent (LoBCoD), operates locally on image patches. Furthermore, we introduce a novel stochastic gradient descent version of LoBCoD for training the convolutional filters. This Stochastic-LoBCoD leverages the benefits of online learning, while being applicable even to a single training image. We demonstrate the advantages of the proposed algorithms for image inpainting and multi-focus image fusion, achieving state-of-the-art results.*

## 1. INTRODUCTION

Sparse representation has been shown to be a very powerful model for many real-world signals, leading to impressive results in various restoration tasks such as denoising [10], deblurring [7], inpainting [11, 25], super-resolution [7, 40] and recognition [37], to name a few. The core assumption of this model is that signals can be expressed as a linear combination of a few columns, also called atoms, taken from a matrix  $\mathbf{D} \in \mathbb{R}^{N \times M}$  termed a dictionary. Concretely, for a signal  $X \in \mathbb{R}^N$ , the model assumption is that  $X = \mathbf{D}\Gamma + V$ , where  $V$  is a noise vector with bounded energy  $\|V\|_2 < \epsilon$ , which allows for a slight deviation from the model and/or may account for noise in the signal. The

vector  $\Gamma \in \mathbb{R}^M$  is the sparse representation of the signal, obtained by solving the pursuit problem [1, 9]:

$$\hat{\Gamma} = \arg \min_{\Gamma} \|\Gamma\|_0 \text{ s.t. } \|X - \mathbf{D}\Gamma\|_2 < \epsilon, \quad (1)$$

where  $\|\Gamma\|_0$  counts the number of non-zeros in  $\Gamma$ . The solution of problem (1) can be approximated using greedy algorithms such as Orthogonal Matching Pursuit (OMP) [4] or convex relaxation algorithms such as Basis Pursuit (BP) [5]. Over the years, various methods have been proposed to adaptively learn the dictionary  $\mathbf{D}$  from real data. Prime examples are K-SVD [1], MOD [12], Double sparsity [29], Online dictionary learning [24], Trainlets [33], and more.

When dealing with high-dimensional signals, learning the dictionary suffers from the curse of dimensionality, rendering this task infeasible. To cope with this problem, many algorithms suggest training a local model on fully-overlapping patches taken from the signal  $X$ . This patch-based technique has gained much popularity due to its simplicity and high-performance [7, 10, 25, 40]. Yet, patch-based approaches are known to be sub-optimal as they ignore the relations between neighboring patches [28, 32].

An alternative approach to meet this challenge is posed by the Convolutional Sparse Coding (CSC) model. This model assumes that the signal can be represented as a superposition of a few local filters, convolved with sparse feature-maps. The CSC model utilizes a structured dictionary (union of narrowly banded convolutional matrices) that facilitates a global handling of the signal. This model has been the subject of an extensive research in the past several years, shown to lead to superior performance in applications such as super-resolution [17], inpainting [18], image separation [26], source separation [20] and audio processing [16].

Contemporary CSC based algorithms often rely on the ADMM [2] formulation for representation-extraction and filter-training. While the majority of works employ ADMM in the Fourier domain [3, 18, 35], a recent approach (SBDL) proposed by Pappyan et al. [26], adopts a local point of view and trains the filters in terms of only local computations in the signal domain. The SBDL algorithm demonstrates

state-of-the-art performance compared to the Fourier-based methods, albeit still relying on the ADMM algorithm. As such, this approach incurs additional memory and sensitivity to additional parameters, it only accommodates a batch-learning mode, and its convergence is questionable<sup>1</sup>.

In this work we propose intuitive and easy-to-implement algorithms, based on the block coordinate descent approach, for solving the global pursuit and the CSC filter learning problems, all done with local computations in the original domain. The proposed pursuit algorithm operates without auxiliary variables nor extra parameters for tuning. We call this algorithm Local Block Coordinate Descent (LoBCoD). In addition, we introduce a stochastic gradient descent variant of LoBCoD for training the convolutional filters. This algorithm leverages the benefits of online learning, while being applicable even to a single training-image. The LoBCoD algorithm and its stochastic version show faster convergence and achieve a better solution to the CSC problem compared to the previous ADMM-based methods (global or local).

The rest of this paper is organized as follows: Section 2 reviews the CSC model and discusses previous methods. The proposed pursuit algorithm is presented in Section 3. In Section 4 we discuss dictionary update methods and introduce the stochastic LoBCoD algorithm. We compare these methods with previously published approaches in section 5. Section 6 extends our methods to image inpainting and multi-focus image fusion, followed by empirical results in Section 7. Section 8 concludes this work.

## 2. Convolutional sparse coding

The CSC model assumes that a signal<sup>2</sup>  $X \in \mathbb{R}^N$  can be represented by the sum of  $m$  convolutions. These are built by feature maps  $\{Z_i\}_{i=1}^m$ , each of length of the original signal  $N$ , convolved with  $m$  small support filters  $\{d_i\}_{i=1}^m$  of length  $n \ll N$ . In the dictionary learning problem, one minimizes the following cost function over both the filters and the feature maps<sup>3</sup>:

$$\min_{d_i, Z_i} \frac{1}{2} \|X - \sum_{i=1}^m d_i * Z_i\|_2^2 + \lambda \sum_{i=1}^m \|Z_i\|_1. \quad (2)$$

Given the filters, the above problem becomes the CSC pursuit task of finding the representations  $\{Z_i\}_{i=1}^m$ . Consider a global dictionary  $\mathbf{D}$  to be the concatenation of  $m$  banded circulant matrices, where each matrix represents a convolution with one filter  $d_i$ . By permuting its columns, the global dictionary  $\mathbf{D}$  consists of all shifted versions of a local dictionary  $\mathbf{D}_L$  of size  $n \times m$ , containing the filters  $\{d_i\}_{i=1}^m$

<sup>1</sup>While SBDL's pursuit method is provably converging, this is no longer true when the dictionary is updated within the ADMM.

<sup>2</sup>The description given focuses on 1D signals for simplicity of the presentation, and all our treatment applies to higher dimensions just as well.

<sup>3</sup>We assume that the filters are normalized to a unit  $l_2$ -norm.

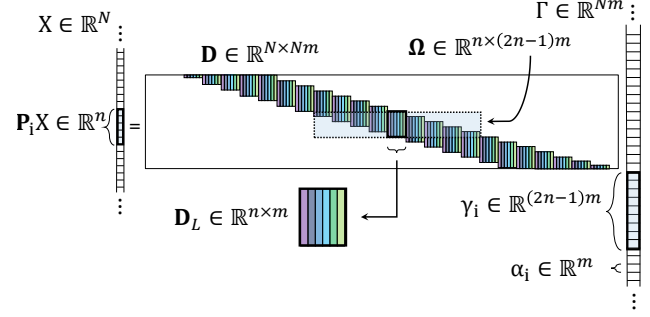


Figure 1: The CSC model and its local components.

as its columns, and the global sparse vector  $\Gamma$  is simply the interlaced concatenation of all the feature maps  $\{Z_i\}_{i=1}^m$ . Such a structure is depicted in Figure 1. Using the above formulation, the convolutional dictionary learning problem (2) can be rewritten as

$$\min_{\mathbf{D}, \Gamma} \frac{1}{2} \|X - \mathbf{D}\Gamma\|_2^2 + \lambda \|\Gamma\|_1. \quad (3)$$

Similar to our earlier comment, when  $\mathbf{D}$  is known, we obtain the CSC pursuit problem, defined as

$$\min_{\Gamma} \frac{1}{2} \|X - \mathbf{D}\Gamma\|_2^2 + \lambda \|\Gamma\|_1. \quad (4)$$

Herein, we review some of the definitions from [27] as they will serve us later for the description of our algorithms.

The global sparse vector  $\Gamma$  can be broken into  $N$  non-overlapping  $m$  dimensional local vectors  $\alpha_i$ , referred to as *needles*. This way, one can express the global vector  $X$  as  $X = \sum_{i=1}^N \mathbf{P}_i^T \mathbf{D}_L \alpha_i$ , where  $\mathbf{P}_i^T \in \mathbb{R}^{N \times n}$  is the operator that positions  $\mathbf{D}_L \alpha_i$  in the  $i$ -th location and pads the rest of the entries with zeros. On the other hand, a patch  $\mathbf{P}_i X = \mathbf{P}_i \mathbf{D} \Gamma$  taken from the signal  $X$  equals to  $\Omega \gamma_i$  (see Figure 1), where  $\Omega \in \mathbb{R}^{n \times (2n-1)m}$  is a *stripe* dictionary containing  $\mathbf{D}_L$  in its center, and  $\gamma_i$  is the *stripe* vector containing the local vector  $\alpha_i$  in its center. In other words, a stripe  $\gamma_i$  is the sparse vector that codes all the content in the patch  $\mathbf{P}_i X$ , whereas a needle  $\alpha_i$  only codes part of the information within it.

The theoretical work in [27] suggested an analysis of the CSC global model, augmented by a localized sparsity measure. Inspired by this analysis, herein we maintain such a local-global decomposition and propose a global algorithm that operates locally on image patches.

## 3. Proposed Method: CSC Pursuit

### 3.1. Local Block Coordinate Descent

In this section we focus on the pursuit of the representations, leaving the study of updating the dictionary for Section 4. The convolutional sparse coding problem presented in the previous section is solved by minimizing the global

objective of Equation (4). In this paper, we adopt a local strategy and split the global sparse vector  $\Gamma$  into local vectors, needles  $\alpha_i$ , and express the global CSC problem in term of such needles and the local dictionary  $\mathbf{D}_L$  by

$$\min_{\mathbf{D}_L, \{\alpha_i\}} \frac{1}{2} \|X - \sum_{i=1}^N \mathbf{P}_i^T \mathbf{D}_L \alpha_i\|_2^2 + \lambda \sum_{i=1}^N \|\alpha_i\|_1. \quad (5)$$

However, rather than optimizing with respect to all the needles together, we suggest to treat the needles sequentially, and optimize with respect to each block  $\alpha_i$  separately. As such, the update rule of each needle can be written as

$$\min_{\alpha_i} \frac{1}{2} \left\| \left( X - \sum_{\substack{j=1 \\ j \neq i}}^N \mathbf{P}_j^T \mathbf{D}_L \alpha_j \right) - \mathbf{P}_i^T \mathbf{D}_L \alpha_i \right\|_2^2 + \lambda \|\alpha_i\|_1. \quad (6)$$

By defining  $R_i = (X - \sum_{\substack{j=1 \\ j \neq i}}^N \mathbf{P}_j^T \mathbf{D}_L \alpha_j)$  as the residual image without the contribution of the needle  $\alpha_i$ , we can rewrite Equation (6) as

$$\min_{\alpha_i} \frac{1}{2} \|R_i - \mathbf{P}_i^T \mathbf{D}_L \alpha_i\|_2^2 + \lambda \|\alpha_i\|_1. \quad (7)$$

While the above minimization involves global variables, such as the residual  $R_i$ , one can show<sup>4</sup> that this can be decomposed into an equivalent and local problem:

$$\min_{\alpha_i} \frac{1}{2} \|\mathbf{P}_i R_i - \mathbf{D}_L \alpha_i\|_2^2 + \lambda \|\alpha_i\|_1. \quad (8)$$

This follows from the observation that the update rule of the needle  $\alpha_i$  is affected only by pixels of the corresponding patch  $\mathbf{P}_i R_i$  (the part that fully overlaps with  $\mathbf{D}_L \alpha_i$ ).

The main idea of the block coordinate descent algorithm is that every step minimizes the overall penalty w.r.t. a certain block of coordinates, while the other ones are set to their most updated values. Following this idea, every local pursuit (8) proceeds by updating the global reconstructed signal  $\hat{X}$  and the global residual  $R = X - \hat{X}$ , as a preprocessing stage that precedes the update of the next needle, based on the most updated values of the previous needles.

An important insight is that needles that have no footprint overlap in the image can be updated efficiently in parallel in the above algorithm without changing the algorithm's outcome. This enables employing efficient batch-implementations of the LARS algorithm [8]. Alternatively, the calculation can be distributed across multiple processors to gain a significant speedup in performance. To formalize these observations, we define the *layer*  $L_i$  as the set of needles that have no induced overlap in the image. We sweep through these layers and update their respective needles in parallel, followed by updating the global reconstructed signal  $\hat{X}$  and the global residual  $R$ . This way, the number of

<sup>4</sup>The proof is provided in the supplementary material.

---

**Algorithm 1:** The stochastic LoBCoD pursuit and dictionary learning algorithm

---

**Input:** signal  $X$ , initial  $\mathbf{D}_L$ , initial needles  $\{\alpha_i^0\}_{i=1}^N$

**Output:** needles  $\{\alpha_i\}_{i=1}^N$ , the trained dictionary  $\mathbf{D}_L$

**Initialization:**  $R = X - \hat{X}$ ,  $\hat{X} = \sum_{i=1}^N \mathbf{P}_i^T \mathbf{D}_L \alpha_i^0$ ,  $k = 0$

**while not converged do**

$k = k + 1$

**for**  $j = 1:n$  **do**

        Computation of the residual:

$$R_j = R + \sum_{i \in L_j} \mathbf{P}_i^T \mathbf{D}_L \alpha_i^{k-1}$$

        Sparse pursuit:  $\forall i \in L_j$  (in parallel)

$$\alpha_i^k = \arg \min_{\alpha_i} \frac{1}{2} \|\mathbf{P}_i R_j - \mathbf{D}_L \alpha_i\|_2^2 + \lambda \|\alpha_i\|_1$$

        Computation of the reconstructed signal:

$$\hat{X} = \hat{X} + \sum_{i \in L_j} \mathbf{P}_i^T \mathbf{D}_L (\alpha_i^k - \alpha_i^{k-1})$$

        Computation of the residual signal:

$$R = X - \hat{X}$$

        Computation of the gradient w.r.t  $\mathbf{D}_L$ :

$$\nabla_{\mathbf{D}_L} = - \sum_{i \in L_j} \mathbf{P}_i R (\alpha_i^k)^T$$

        Dictionary update:

$$\mathbf{D}_L = \mathcal{P}_{\mathbb{L}}[\mathbf{D}_L - \eta \nabla_{\mathbf{D}_L}]$$

**end**

**end**

---

the layers imposes the number of the inner iterations, which will determine the complexity of our final algorithm. In that manner, the number of the inner iterations depends only on the patch size; for  $\sqrt{n} \times \sqrt{n}$  patches, the number of layers is  $n$ . This pursuit algorithm is presented in Algorithm 1. Note that this algorithm can clearly be extended to iterate over multiple signals, but for the sake of brevity we assume that the data corresponds to an individual signal  $X$ .

### 3.2. Boundary Conditions and Initialization

In the formulation of the CSC model, as shown in Figure 1, we assumed that the dictionary is comprised of a set of banded circulant matrices, which impose a circulant boundary conditions on the signals. In practice, however, signals and images do not exhibit circulant boundary behavior. Therefore, our model incorporates a preemptive treatment of the boundaries. We adopt a similar approach to [26], in which the signal boundaries are padded with  $n - 1$  elements prior to decomposing it with the model. At the end of the process, we discard the added padding by cropping the  $n - 1$  boundary elements from the reconstructed signal and from the resulting feature maps (sparse representation).

Another beneficial preprocess step is needles initialization. A good initialization would equally spread the contribution of the needles towards signal reconstruction. With that goal, we set the initial value of each needle  $\alpha_i$  to be the sparse representation of  $\frac{1}{n}\mathbf{P}_i X$ , i.e its relative portion of the corresponding patch. This can be done by solving the following local pursuit for every needle:

$$\alpha_i^0 = \arg \min_{\alpha_i} \frac{1}{2} \left\| \frac{1}{n} \mathbf{P}_i X - \mathbf{D}_L \alpha_i \right\|_2^2 + \lambda \|\alpha_i\|_1, \quad (9)$$

as a preprocess stage of our algorithm.

## 4. CSC Dictionary Learning

When addressing the question of learning the CSC filters, the common strategy is to alternate between sparse-coding and dictionary update steps for a fixed number of iterations. The dictionary update step aims to find the minimum of the quadratic term of Equation (5) subject to the constraint of normalized dictionary columns:

$$\begin{aligned} \min_{\mathbf{D}_L} \quad & \frac{1}{2} \left\| X - \sum_{i=1}^N \mathbf{P}_i^T \mathbf{D}_L \alpha_i \right\|_2^2 \\ \text{s.t.} \quad & \{\|d_i\|_2 = 1\}_{i=1}^m. \end{aligned} \quad (10)$$

One can do so in a batch manner which requires access to the entire dataset at every iteration, or in an online (stochastic) manner that enables access to only small part of the dataset at every update step. This way it is also applicable for streaming data scenarios, when the probability distribution of the data changes over time.

### 4.1. Batch Update

Usually, for offline applications where the whole dataset is given and can be stored in memory, the batch approach is generally simpler, and thus we start with its description. The typical approach is to alternate between sparse coding (4) and dictionary update (10). For the latter, solving problem (10) requires finding the optimum  $\mathbf{D}_L$  that satisfies the normalization constraint. One can find this optimal solution using projected steepest descent: perform steepest descent with a small step size and project the solution to the constraint set after each iteration, until convergence. To that end, the gradient of the objective in problem (10) w.r.t.  $\mathbf{D}_L$  is<sup>5</sup>:

$$\nabla_{\mathbf{D}_L} = - \sum_{i=1}^N \mathbf{P}_i (X - \hat{X}) \cdot \alpha_i^T. \quad (11)$$

The final update step for the local dictionary  $\mathbf{D}_L$  is obtained by advancing in the direction of this gradient (11) and normalizing the columns of the resulting  $\mathbf{D}_L$  in each iteration, until convergence.

<sup>5</sup>This derivation can be found in the supplementary material.

This batch dictionary update rule follows the line of thought of the MOD algorithm [12], and thus improves the solution in each step. However, it exhibits a very slow convergence rate since each dictionary update can be performed only after finishing the entire sparse coding (pursuit) stage, which is markedly inefficient, as the pursuit is the most time consuming part of the algorithm. This brings us to the Stochastic-LoBCoD alternative.

### 4.2. Local Stochastic Gradient Descent Approach

The traditional Stochastic Gradient Descent (SGD) approach restricts the computation of the gradient to a subset of the data and advances in the direction of this noisy gradient with every update step. Building upon this concept and the fact that Equation (11) reveals a separable gradient w.r.t the patches and their corresponding needles, we can update the dictionary in a stochastic manner. Rather than concluding the entire pursuit stage and then advancing in the direction of the global gradient, we can take a small step size  $\eta$  and update the dictionary after finding the sparse representation of only a small group of needles. According to Section 3, every iteration updates a group of needles, referred to as a layer  $L_i$ , which in turn could now serve to update the dictionary. This way, our algorithm converges faster and adopts the stochastic behavior of the SGD while still operating on a single image.

The filters should be normalized after every dictionary update by projecting them onto the  $l_2$  unit ball. Here, due to the choice of small step size, we simply normalize the atoms after every dictionary update:

$$\mathbf{D}_L = \mathcal{P}_1[\mathbf{D}_L - \eta \nabla_{\mathbf{D}_L}].$$

Where  $\mathcal{P}_1[\cdot]$  denotes the operator that projects the dictionary atoms onto the unit ball. The final algorithm that incorporates the dictionary update is summarized in Algorithm 1.

Note that, although this dictionary update rule introduces an extra parameter (the step size  $\eta$ ), determining its value is rather intuitive and can be performed automatically by setting it to  $1 - 2\%$  of the norm of the gradient. Furthermore, this update rule may also leverage any stochastic optimization algorithm such as Momentum, Adagrad, Adadelta, Adam [30] etc., with their authors' recommended parameter values. This choice of parameter setting is sufficient, as will be demonstrated empirically in Section 7. In the rest of this work we will use this dictionary update rule, as it shows superior results.

## 5. Relation to Other Methods

In this section we describe the advantages of the proposed approach over Fourier and ADMM based methods.

**Parallel computation:** Our algorithm is trivial to parallelize efficiently across multiple processors by virtue of operating directly on the image patches. One can split the

Method	Time Complexity
[22] (Sparse)	$\underbrace{qImN\log(N)}_{\text{T-pursuit}} + \underbrace{INmk}_{\text{SGD using sparse matrix}}$
[22] (Freq.)	$\underbrace{qImN\log(N)}_{\text{T-pursuit}} + \underbrace{ImN\log(N) + INm}_{\text{SGD in the Fourier domain}}$
SBDL	$\underbrace{INnm + IN(k^3 + mk^2)}_{\text{LARS}} + \underbrace{nm^2}_{\text{Gram}} + \underbrace{INk(n+m) + nm^2}_{\text{K-SVD}}$
Ours	$\underbrace{INnm + IN(k^3 + mk^2)}_{\text{LARS}} + \underbrace{n^2m^2}_{\text{Gram}} + \underbrace{IN(n + nk + m)}_{\text{Stochastic-LoBCoD}}$

Table 1: Complexity analysis of our method compared to the online algorithms presented in [22] and the SBDL algorithm [26]. I: number of signals, N: signal dimension, m: number of filters, n: patch size, k: maximum number of non-zeros per needle, q: number of inner iterations for the pursuit algorithm. The dominant terms are in red color.

computations between  $N/n$  processors, corresponding to the number of the needles in every layer.

**Online learning:** The proposed algorithm, due to its local stochastic manner, can work in a streaming mode, where the probability distribution of the patches varies over time. Another aspect of this advantage is our ability to run in an online manner, even for a single input image. This stands in sharp contrast to other recent online methods [21, 34] that allow for online training but only in the case of streaming images. Other approaches took a step further and proposed partitioning the image into sub-images [22], but this is still far from our approach, which can stochastically estimate the gradient for each needle.

**Parameter free:** Contrary to ADMM-based approaches, our algorithm is unhindered by cumbersome manual parameter-tuning at the pursuit stage. Moreover, it benefits from an intuitively tuned parameter (the step size  $\eta$ ) in the dictionary learning stage, as described as Section 4.

**Memory efficient:** Our algorithm has better storage complexity compare to the ADMM-based approaches [18, 26] since the update of the sparse vector is performed in-place and does not require any auxiliary variables.

Table 1 compares the complexity<sup>6</sup> of executing an epoch in our algorithm with that of the batch algorithm in [26] and the online algorithms in [22]. The conclusion is that our algorithm scales linearly with the global dimension N, while the competing online algorithms grow as  $O(N\log(N))$ .

## 6. Image Processing via CSC

Having established the foundations for our algorithms, we now detail their extended variants for tackling the task of image inpainting and multi-focus image fusion. We also present adaptations of our algorithm for tackling the tasks of multi-exposure image fusion and salt-and-pepper text image denoising in the supplementary material.

<sup>6</sup>The full explanation can be found in the supplementary material.

### 6.1. Image Inpainting

The task of image inpainting pertains to filling-in missing pixels at known locations in the image. Assume we are given a corrupted image  $Y = \mathbf{A}X$ , where  $\mathbf{A} \in \mathbb{R}^{N \times N}$  is a binary diagonal matrix that represents the degradation operator, so that  $\mathbf{A}(i, i) = 0$  implies that the pixel  $x_i$  is masked. The goal of image inpainting is to reconstruct the original image  $X$ . Using the CSC formulation, this can be performed by first solving the following problem:

$$\min_{\Gamma} \frac{1}{2} \|Y - \mathbf{A}\mathbf{D}\Gamma\|_2^2 + \lambda \|\Gamma\|_1, \quad (12)$$

and then taking the found representation  $\Gamma$  and multiplying by  $\mathbf{D}$ . By applying the steps described in Section 3, we split the above global optimization problem into a series of more manageable problems, each acting on a block of coordinates, i.e. a *needle*. This yields the following version of Equation (8):

$$\min_{\alpha_i} \frac{1}{2} \|\mathbf{P}_i R_i - \mathbf{A}_i \mathbf{D}_L \alpha_i\|_2^2 + \lambda \|\alpha_i\|_1. \quad (13)$$

Here,  $\mathbf{A}_i = \mathbf{P}_i \mathbf{A} \mathbf{P}_i^T$  is the operator that masks the corresponding  $i$ -th patch, and  $R_i = (Y - \mathbf{A} \sum_{j=1, j \neq i}^N \mathbf{P}_j^T \mathbf{D}_L \alpha_j)$  is the residual between the corrupted image and the degraded version of the reconstructed image, where the residual  $R_i$  does not account for the needle  $\alpha_i$ . As mentioned in Section 3, we parallelize the computations of the needles that comprised each layer. The dictionary  $\mathbf{D}_L$  can be pretrained on an external, uncorrupted dataset or trained on the corrupted image directly using the following gradient:

$$\nabla_{\mathbf{D}_L} = - \sum_{i \in L_j} \mathbf{P}_i \mathbf{A}^T (Y - \mathbf{A} \hat{X}) \cdot \alpha_i^T, \quad (14)$$

where  $\hat{X} = \sum_{j=1}^N \mathbf{P}_j^T \mathbf{D}_L \alpha_j$  is the reconstructed image. The gradient derivation above is identical to that in Section 4, but with the exception of incorporating the mask  $\mathbf{A}$ .

### 6.2. Multi-focus image fusion

Image fusion techniques aim to integrate complimentary information from multiple images, captured with different focal settings, into an all-in-focus image of higher quality. Many patch-based sparse formulations were proposed to address this task, such as choose-max OMP [38], simultaneous OMP [39], and coupled sparse representation [14]. In this work, we adopt a similar scheme to [23], which utilizes the CSC for tackling the task of image-fusion, but with the distinction of solving a unified minimization problem.

Assume we are given a set of source images  $\{Y^k\}_{k=1}^L$  to fuse, as well as a set of pretrained filters  $\{d_i\}_{i=1}^m$ . We start by decomposing each image  $Y^k$  into a base component  $Y_b^k$  and an edge component  $Y_e^k$  by imposing distinctive priors.



For the base component  $Y_b^k$  we penalize the  $l_2$  norm of its gradient, while for the edge component we employ the CSC model such that  $Y_e^k = \sum_{i=1}^m d_i * Z_i^k$ . Practically, for each image  $Y^k$  we solve the unified minimization problem:

$$\min_{\{Z_i^k\}, Y_b^k} \frac{1}{2} \|Y^k - \sum_{i=1}^m d_i * Z_i^k - Y_b^k\|_2^2 + \lambda \sum_{i=1}^m \|Z_i^k\|_1 + \mu \frac{1}{2} \|\nabla Y_b^k\|_2^2. \quad (15)$$

This is done by alternating between minimizing w.r.t. the base component  $Y_b^k$  and the feature maps  $\{Z_i^k\}_{i=1}^m$ . The former boils down to a least square problem, and the latter is solved using our LoBCoD algorithm<sup>7</sup>.

After decomposing all the images, we aim to fuse their components. For each image, we build an activity map  $\tilde{A}^k$ , as the sum of the absolute values of  $\{Z_i^k\}_{i=1}^m$ . For robustness, we convolve  $\tilde{A}^k$  with a uniform kernel  $U_s \in \mathbb{R}^{s \times s}$ :

$$\tilde{A}^k(u, v) = \sum_{i=1}^m \|Z_i^k(u, v)\|_1, \quad A^k = \tilde{A}^k * U_s. \quad (16)$$

Based on the observation that a significant value in the activity map  $A^k$  indicates a sharp region in the corresponding image  $Y^k$ , we reconstruct the all-in-focus components by assembling the most prominent regions based on their values in the corresponding activity maps:

$$Z_i^f(u, v) = Z_i^{k^*}(u, v), \quad Y_b^f(u, v) = Y_b^{k^*}(u, v), \quad (17)$$

$$k^* = \arg \max_k (A^k(u, v)).$$

where  $\{Z_i^f\}_{i=1}^m$  and  $Y_b^f$  are the feature maps and the base component of the fused image  $Y^f$ . Finally, the fusion result is obtained by gathering its components:

$$Y^f = Y_b^f + \sum_{i=1}^m d_i * Z_i^f. \quad (18)$$

## 7. Experiments

The full LoBCoD implementation, documentation and demos that reproduce our results, are available online<sup>8</sup>.

### 7.1. Run Time Comparison

To begin with, and to provide a comparison to other state of the art methods, we evaluate the performance of the proposed algorithm for solving Equation (5) against state of the art batch algorithms for CSC: the SBDL algorithm [26], the algorithm in [36] and the algorithm presented in [15], all using the same settings on the Fruit dataset [13]. For learning the dictionary, we used the ADAM and the Momentum

algorithms<sup>9</sup> [30]. Figure 2 presents a comparison of the objective (2) as a function of time for each of the competing algorithms, showing that our method achieves the fastest convergence. Figure 4 shows the obtained dictionaries.

We also compared our method to the online stochastic gradient descent (SGD) based algorithms in [22], which operate in the spatial and in the Fourier domains. Here we randomly selected a training set of 40 images, and a test set of 5 different images from the MIRFLICKER-1M dataset [19]. Figure 3 presents the objective of the test set as a function of time, showing that our algorithm converges faster. Figure 5 shows the dictionaries obtained by the three methods, illustrating similar quality. Note that our algorithm is capable of operating online even if trained on one image, a possibility that is not supported by [22].

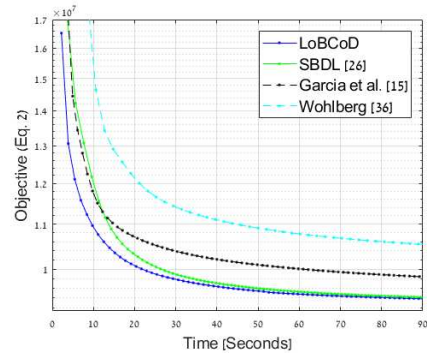


Figure 2: Run time comparison between our method and the batch methods in [26], [15] and [36].

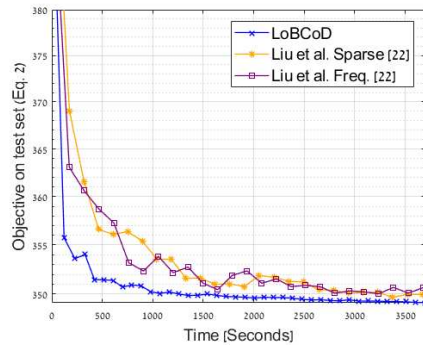


Figure 3: Run time comparison between our method and the online algorithms in [22].

### 7.2. Image Inpainting

We apply our algorithm to the task of image inpainting, as described in Section 6.1, and compare our results to [26]. In Section 6.1 we described two viable methods for training the dictionary; utilizing an external dataset or training directly on the corrupted source image. For the former, we used the Fruit dataset [13] for both algorithms, as shown in

<sup>7</sup>Additional information can be found in the supplementary material.

<sup>8</sup><https://github.com/EvZissel/LoBCoD>

<sup>9</sup>The parameter settings are described in the supplementary material.

Table 2: Inpainting comparison [dB] between the proposed LoBCoD and the SBDL [26] algorithms.

	Barbara	Boat	House	Lena	Peppers	C.man	Couple	Finger	Hill	Man	Montage
SBDL external	30.41	31.76	36.17	35.92	<b>33.69</b>	28.76	32.16	30.91	33.12	33.04	28.93
Proposed external	<b>30.93</b>	<b>31.82</b>	<b>36.58</b>	<b>36.15</b>	33.54	<b>28.88</b>	<b>32.46</b>	<b>31.75</b>	<b>33.25</b>	<b>33.18</b>	<b>29.18</b>
SBDL internal	31.98	32.04	36.19	36.01	34.03	28.85	32.18	30.96	33.21	32.99	28.95
Proposed internal	<b>32.50</b>	<b>32.27</b>	<b>36.74</b>	<b>36.17</b>	<b>34.48</b>	<b>29.04</b>	<b>32.56</b>	<b>31.76</b>	<b>33.42</b>	<b>33.25</b>	<b>29.23</b>



(a) Ours (b) SBDL [26] (c) [15]

Figure 4: Comparison between the dictionaries obtained using the Stochastic-LoBCoD method vs. the methods in [26] and [15] on the Fruit dataset [13].



(a) Ours (b) [22] Sparse (c) [22] Freq.

Figure 5: Comparison between the dictionaries obtained using the Stochastic-LoBCoD method vs. the online methods in [22] on the MIRFLICKR-1M dataset [19].

Figure 4. The corrupted images were created by applying a randomly generated mask with 50% missing pixels on the original images. All the corrupted test images were mean-subtracted prior to applying both algorithms by subtracting the patch-average of the unmasked pixels. In addition, we tuned  $\lambda$  in Equation (12) for every corrupted test image, to account for their varying complexity. The top two rows of Table 2 present the results using an external dataset in terms of peak signal-to-noise ratio (PSNR) on a set of 11 test images, showing that our method leads to better results. Next, we train the dictionary of both algorithms on the corrupted image itself. The results are presented at the bottom two rows of Table 2, indicating that the Stochastic-LoBCoD algorithm achieves better results.

### 7.3. Multi-focus image fusion

We conclude by applying our LoBCoD algorithm to the task of multi-focus image fusion, as described in the previous section. We evaluate our proposed method using synthetic data, as well as data from a real dataset, and compare our results to [23]. The dictionaries of both methods were pretrained on the Fruit dataset [13].

For the synthetic experiment, we extracted a portion of



(a) Barbara in-focus (b) Barbara out-of-focus



(c) [23] 41.96dB (d) Proposed 42.27dB



(e) Butterfly in-focus (f) Butterfly out-of-focus



(g) [23] 35.30dB (h) Proposed 36.01dB

Figure 6: Fusion performance comparison on synthetic images. The PSNR values were computed between the reconstructed and the original images.

the standard image Barbara and created two input images of blurred foreground and blurred background. Image blurring was performed using a  $9 \times 9$  Gaussian blur kernel with  $\sigma = 2$ . We repeated the same procedure on the image Butterfly<sup>10</sup>, using a  $16 \times 16$  Gaussian blur kernel with  $\sigma = 4$ .

<sup>10</sup>The image was taken from the dataset in [6].

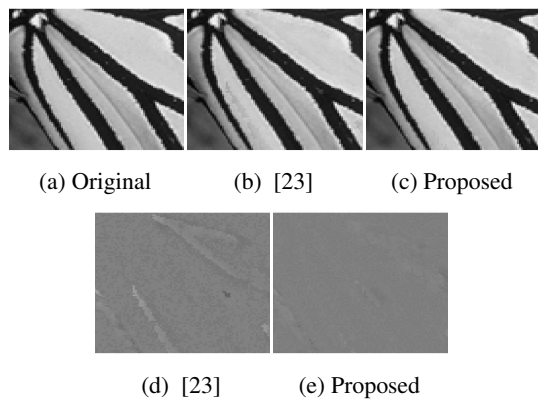


Figure 7: Zoom in on the fusion results of the image Butterfly. Figures (d) and (e) present the error images of (b) and (c), respectively. The error images were computed between the fusion results and the original image (a).

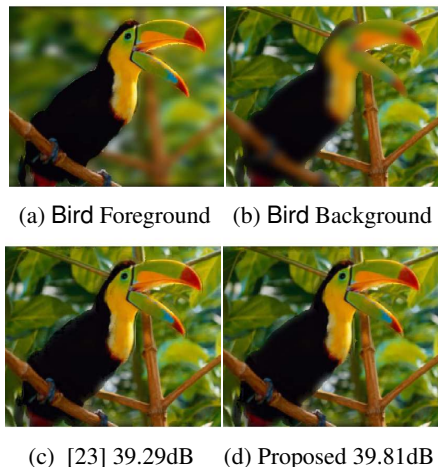


Figure 8: Fusion comparison between the proposed method and the method in [23] on the image Bird.

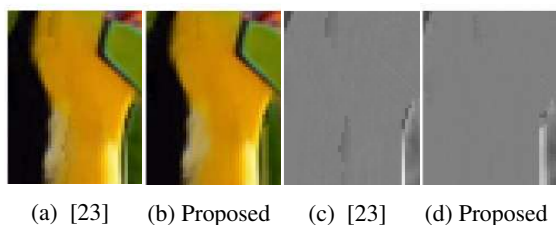


Figure 9: Zoom in on the fusion results of the image Bird. Figure (c) and (d) present the L channel error compared to that of the original image.

Both sets of synthetic blurred images are presented in Figure 6, alongside their reconstructed images, and the corresponding PSNR values. The resulting images demonstrate that our approach leads to visually and quantitatively better results. Figure 7 presents a zoom-in view of our reconstructed image Butterfly, compared to the result of [23] and the original image; showing that for images with prominent blur our method achieves visually better results.

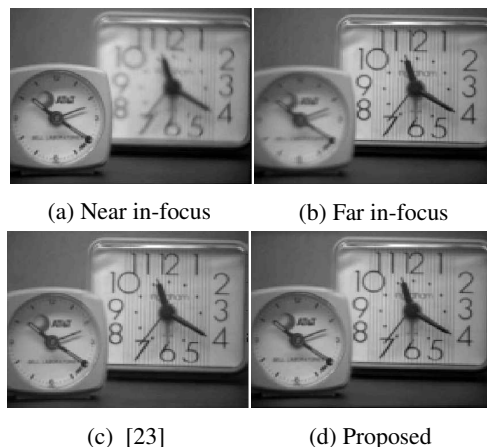


Figure 10: Fusion examples of real images taken from [31].

We adapted our approach for fusion of colored images, and showed the result on the image Bird<sup>11</sup>. We generate a pair of blurred background and foreground color images by applying a  $16 \times 16$  Gaussian blur kernel with  $\sigma = 4$  on each RGB channel. We chose to blur the image in the RGB color space to emulate a blur of a camera. Afterwards, both blurred colored images were treated by transforming them to the Lab color space. The PSNR for the Bird image was computed between the L channels of the original and the reconstructed images. We present the results together with their PSNR values in Figure 8 and 9, showing that our approach leads to visually and quantitatively better results.

Lastly, for the real dataset experiment, we ran our proposed algorithm on the image-pair Clocks taken from [31]. Figure 10 presents the resulting fused images, showing comparable results on this image-pair.

## 8. Conclusions

In this work we have introduced the local block coordinate descent (LoBCoD) algorithm for performing pursuit for the global CSC model, while operating locally on image patches. We demonstrated its advantages over contending state-of-the-art methods in terms of memory requirements, efficient parallel computation, and its exemption from meticulous manual tuning of parameters. In addition, we proposed a stochastic gradient descent version (Stochastic-LoBCoD) of this algorithm for training the convolutional filters. We highlighted its unique qualities as an online algorithm that retains the ability to act on a single image. Finally, we illustrated the advantages of the proposed algorithm on a set of applications and compared it with competing state-of-the-art methods.

## 9. Acknowledgments

This work was supported by the Israel Science Foundation (ISF) under Grant 335/18.

<sup>11</sup>The image was taken from the dataset in [7]



## References

- [1] M. Aharon, M. Elad, and A. Bruckstein. K-svd: An algorithm for designing overcomplete dictionaries for sparse representation. *IEEE Transactions on signal processing*, 54(11):4311–4322, 2006.
- [2] S. Boyd, N. Parikh, E. Chu, B. Peleato, J. Eckstein, et al. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends® in Machine learning*, 3(1):1–122, 2011.
- [3] H. Bristow, A. Eriksson, and S. Lucey. Fast convolutional sparse coding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 391–398, 2013.
- [4] S. Chen, S. A. Billings, and W. Luo. Orthogonal least squares methods and their application to non-linear system identification. *International Journal of control*, 50(5):1873–1896, 1989.
- [5] S. S. Chen, D. L. Donoho, and M. A. Saunders. Atomic decomposition by basis pursuit. *SIAM review*, 43(1):129–159, 2001.
- [6] W. Dong, L. Zhang, G. Shi, and X. Li. Nonlocally centralized sparse representation for image restoration. *IEEE Transactions on Image Processing*, 22(4):1620–1630, 2013.
- [7] W. Dong, L. Zhang, G. Shi, and X. Wu. Image deblurring and super-resolution by adaptive sparse domain selection and adaptive regularization. *IEEE Transactions on Image Processing*, 20(7):1838–1857, 2011.
- [8] B. Efron, T. Hastie, I. Johnstone, R. Tibshirani, et al. Least angle regression. *The Annals of statistics*, 32(2):407–499, 2004.
- [9] M. Elad. *Sparse and Redundant Representations: From Theory to Applications in Signal and Image Processing*. Springer, 2010.
- [10] M. Elad and M. Aharon. Image denoising via sparse and redundant representations over learned dictionaries. *IEEE Transactions on Image processing*, 15(12):3736–3745, 2006.
- [11] M. Elad, J.-L. Starck, P. Querre, and D. L. Donoho. Simultaneous cartoon and texture image inpainting using morphological component analysis (mca). *Applied and Computational Harmonic Analysis*, 19(3):340–358, 2005.
- [12] K. Egan, S. O. Aase, and J. H. Husoy. Method of optimal directions for frame design. In *Acoustics, Speech, and Signal Processing, 1999. Proceedings., 1999 IEEE International Conference on*, volume 5, pages 2443–2446. IEEE, 1999.
- [13] R. Fergus, M. D. Zeiler, G. W. Taylor, and D. Krishnan. Deconvolutional networks. In *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 00, pages 2528–2535, 2010.
- [14] R. Gao and S. A. Vorobyov. Multi-focus image fusion via coupled sparse representation and dictionary learning. *arXiv preprint arXiv:1705.10574*, 2017.
- [15] C. Garcia-Cardona and B. Wohlberg. Subproblem coupling in convolutional dictionary learning. In *Image Processing (ICIP), 2017 IEEE International Conference on*, pages 1697–1701. IEEE, 2017.
- [16] R. Grosse, R. Raina, H. Kwong, and A. Y. Ng. Shift-invariant sparse coding for audio classification. In *The Twenty-Third Conference on Uncertainty in Artificial Intelligence (UAI)*, pages 149–158, 2007.
- [17] S. Gu, W. Zuo, Q. Xie, D. Meng, X. Feng, and L. Zhang. Convolutional sparse coding for image super-resolution. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1823–1831, 2015.
- [18] F. Heide, W. Heidrich, and G. Wetzstein. Fast and flexible convolutional sparse coding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5135–5143, 2015.
- [19] M. J. Huiskes, B. Thomee, and M. S. Lew. New trends and ideas in visual concept detection: the mir flickr retrieval evaluation initiative. In *Proceedings of the international conference on Multimedia information retrieval*, pages 527–536. ACM, 2010.
- [20] H.-W. Liao and L. Su. Monaural source separation using ramanujan subspace dictionaries. *IEEE Signal Processing Letters*, 25(8), 2018.
- [21] J. Liu, C. Garcia-Cardona, B. Wohlberg, and W. Yin. On-line convolutional dictionary learning. In *Image Processing (ICIP), 2017 IEEE International Conference on*, pages 1707–1711. IEEE, 2017.
- [22] J. Liu, C. Garcia-Cardona, B. Wohlberg, and W. Yin. First-and second-order methods for online convolutional dictionary learning. *SIAM Journal on Imaging Sciences*, 11(2):1589–1628, 2018.
- [23] Y. Liu, X. Chen, R. K. Ward, and Z. J. Wang. Image fusion with convolutional sparse representation. *IEEE signal processing letters*, 23(12):1882–1886, 2016.
- [24] J. Mairal, F. Bach, J. Ponce, and G. Sapiro. Online dictionary learning for sparse coding. In *Proceedings of the 26th annual international conference on machine learning*, pages 689–696. ACM, 2009.
- [25] J. Mairal, M. Elad, and G. Sapiro. Sparse representation for color image restoration. *IEEE Transactions on image processing*, 17(1):53–69, 2008.
- [26] V. Papan, Y. Romano, J. Sulam, and M. Elad. Convolutional dictionary learning via local processing. In *ICCV*, pages 5306–5314, 2017.
- [27] V. Papan, J. Sulam, and M. Elad. Working locally thinking globally: Theoretical guarantees for convolutional sparse coding. *IEEE Transactions on Signal Processing*, 65(21):5687–5701, 2017.
- [28] Y. Romano and M. Elad. Patch-disagreement as away to improve k-svd denoising. In *Acoustics, Speech and Signal Processing (ICASSP), 2015 IEEE International Conference on*, pages 1280–1284. IEEE, 2015.
- [29] R. Rubinstein, M. Zibulevsky, and M. Elad. Double sparsity: Learning sparse dictionaries for sparse signal approximation. *IEEE Transactions on signal processing*, 58(3):1553–1564, 2010.
- [30] S. Ruder. An overview of gradient descent optimization algorithms. *arXiv preprint arXiv:1609.04747*, 2016.
- [31] S. Savić. Multifocus image fusion based on empirical mode decomposition. In *Twentieth International Electro technical and Computer Science Conference*, 2011.

- [32] J. Sulam and M. Elad. Expected patch log likelihood with a sparse prior. In *International Workshop on Energy Minimization Methods in Computer Vision and Pattern Recognition*, pages 99–111. Springer, 2015.
- [33] J. Sulam, B. Ophir, M. Zibulevsky, and M. Elad. Trainlets: Dictionary learning in high dimensions. *IEEE Transactions on Signal Processing*, 64(12):3180–3193, 2016.
- [34] Y. Wang, Q. Yao, J. T. Kwok, and L. M. Ni. Scalable online convolutional sparse coding. *IEEE Transactions on Image Processing*, 2018.
- [35] B. Wohlberg. Efficient convolutional sparse coding. In *Acoustics, Speech and Signal Processing (ICASSP), 2014 IEEE International Conference on*, pages 7173–7177. IEEE, 2014.
- [36] B. Wohlberg. Boundary handling for convolutional sparse representations. In *Image Processing (ICIP), 2016 IEEE International Conference on*, pages 1833–1837. IEEE, 2016.
- [37] J. Wright, A. Y. Yang, A. Ganesh, S. S. Sastry, and Y. Ma. Robust face recognition via sparse representation. *IEEE transactions on pattern analysis and machine intelligence*, 31(2):210–227, 2009.
- [38] B. Yang and S. Li. Multifocus image fusion and restoration with sparse representation. *IEEE Transactions on Instrumentation and Measurement*, 59(4):884–892, 2010.
- [39] B. Yang and S. Li. Pixel-level image fusion with simultaneous orthogonal matching pursuit. *Information fusion*, 13(1):10–19, 2012.
- [40] J. Yang, J. Wright, T. S. Huang, and Y. Ma. Image super-resolution via sparse representation. *IEEE transactions on image processing*, 19(11):2861–2873, 2010.

# A Local Block Coordinate Descent Algorithm for the CSC Model

## Supplementary Material

Ev Zisselman

ev\_zis@campus.technion.ac.il

Jeremias Sulam

jsulam1@jhu.edu

Michael Elad

elad@cs.technion.ac.il

This Supplementary material elaborates on the LoBCoD algorithm and our experiments. In Section A we provide proofs for our derivations of the LoBCoD algorithm and its stochastic version. Section B analyzes the complexity of our approach and compares it to contending batch and on-line algorithms. Section C details the algorithm for image fusion and in Section D we specify the implementation details and parameters settings of our experiments. Lastly, in Section E we demonstrate the advantages of LoBCoD on two additional applications: multi-exposure image fusion and salt-and-pepper text image denosing.

### A. Mathematical proofs

#### A.1. Transitioning from a high dimensional problem to a low dimensional problem

**Transition from minimization problem (7) to (8) in the main paper.** We denote  $p_i$  as the patch that fully contains the slice  $s_i = \mathbf{D}_L \alpha_i$ , and define a patch-layer  $L_{\tilde{i}}$  as the set of non-overlapping patches taken from the image that contains the patch  $p_i$ . We can write the identity matrix as a sum of non-overlapping patch-extraction matrices  $\sum_{k \in L_{\tilde{i}}} \mathbf{P}_k^T \mathbf{P}_k = \mathbf{I}$ , where  $\mathbf{P}_i$  is one of these matrices. By using these definitions and writing the definition of  $R_i$  (the residual image without the contribution of the  $i$ -th needle) we can write the  $l_2$  fidelity term of Equation (7) as

$$\begin{aligned} \arg \min_{\alpha_i} \frac{1}{2} \left\| \sum_{k \in L_{\tilde{i}}} \mathbf{P}_k^T \mathbf{P}_k R_i - \mathbf{P}_i^T \mathbf{D}_L \alpha_i \right\|_2^2 &= \\ \arg \min_{\alpha_i} \frac{1}{2} \left\| \sum_{k \in L_{\tilde{i}}} \mathbf{P}_k^T \mathbf{P}_k (X - \sum_{\substack{j=1 \\ j \neq i}}^N \mathbf{P}_j^T \mathbf{D}_L \alpha_j) - \mathbf{P}_i^T \mathbf{D}_L \alpha_i \right\|_2^2 &= \\ \arg \min_{\alpha_i} \frac{1}{2} \left\| \sum_{\substack{k \in L_{\tilde{i}} \\ k \neq i}} \mathbf{P}_k^T \mathbf{P}_k (X - \sum_{\substack{j=1 \\ j \neq i}}^N \mathbf{P}_j^T \mathbf{D}_L \alpha_j) + \right. \\ \left. \mathbf{P}_i^T (\mathbf{P}_i (X - \sum_{\substack{j=1 \\ j \neq i}}^N \mathbf{P}_j^T \mathbf{D}_L \alpha_j) - \mathbf{D}_L \alpha_i) \right\|_2^2. \end{aligned} \quad (\text{A.1})$$

Since the patch-extraction matrix  $\mathbf{P}_i^T$  is orthogonal to all the matrices  $\mathbf{P}_k^T$  for  $k \neq i$ , the previous problem (A.1) is equal to

$$\begin{aligned} \arg \min_{\alpha_i} \frac{1}{2} \left\| \sum_{\substack{k \in L_{\tilde{i}} \\ k \neq i}} \mathbf{P}_k^T \mathbf{P}_k (X - \sum_{\substack{j=1 \\ j \neq i}}^N \mathbf{P}_j^T \mathbf{D}_L \alpha_j) \right\|_2^2 + \\ \left\| \mathbf{P}_i^T (\mathbf{P}_i (X - \sum_{\substack{j=1 \\ j \neq i}}^N \mathbf{P}_j^T \mathbf{D}_L \alpha_j) - \mathbf{D}_L \alpha_i) \right\|_2^2. \end{aligned} \quad (\text{A.2})$$

Note that the first term of the above objective does not depend on  $\alpha_i$ , and thus we can ignore this term in our minimization of the objective:

$$\arg \min_{\alpha_i} \frac{1}{2} \left\| \mathbf{P}_i^T (\mathbf{P}_i (X - \sum_{\substack{j=1 \\ j \neq i}}^N \mathbf{P}_j^T \mathbf{D}_L \alpha_j) - \mathbf{D}_L \alpha_i) \right\|_2^2. \quad (\text{A.3})$$

In addition, the matrix  $\mathbf{P}_i^T$  translates a patch-size vector to the  $i$ -th position in the global vector padded with zeros. Hence, we can ignore all the zero-entries in the resulting vector, and the problem becomes equivalent to solving the following reduced minimization problem:

$$\begin{aligned} \arg \min_{\alpha_i} \frac{1}{2} \left\| \mathbf{P}_i (X - \sum_{\substack{j=1 \\ j \neq i}}^N \mathbf{P}_j^T \mathbf{D}_L \alpha_j) - \mathbf{D}_L \alpha_i \right\|_2^2 &= \\ \arg \min_{\alpha_i} \frac{1}{2} \left\| \mathbf{P}_i R_i - \mathbf{D}_L \alpha_i \right\|_2^2, \end{aligned} \quad (\text{A.4})$$

where  $R_i = X - \sum_{\substack{j=1 \\ j \neq i}}^N \mathbf{P}_j^T \mathbf{D}_L \alpha_j$ .

#### A.2. Local dictionary update - gradient calculation

**The derivation of the gradient, Equation (11) in the main paper.** We can rewrite  $\mathbf{D}_L$  as a column vector  $d_L$  and write problem (10) as:

$$\arg \min_{\mathbf{D}_L} \frac{1}{2} \left\| X - \sum_{i=1}^N \mathbf{P}_i^T (\alpha_i \otimes \mathbf{I}_{n \times n}) d_L \right\|_2^2, \quad (\text{A.5})$$

where we denote  $\otimes$  as the Kronecker matrix product and apply property no. (40) from [9].

By defining  $\mathbf{A}_i = (\alpha_i \otimes \mathbf{I}_{n \times n})$ , the above problem can be rewritten as

$$\arg \min_{\mathbf{D}_L} \frac{1}{2} \left\| X - \sum_{i=1}^N \mathbf{P}_i^T \mathbf{A}_i d_L \right\|_2^2. \quad (\text{A.6})$$

Now it easy to see that the gradient of problem (A.6) w.r.t.  $d_L$  is given by

$$\begin{aligned} \nabla_{d_L} &= - \left( \sum_{i=1}^N \mathbf{P}_i^T \mathbf{A}_i \right)^T \left( X - \sum_{i=1}^N \mathbf{P}_i^T \mathbf{A}_i d_L \right) = \\ &= - \sum_{i=1}^N \mathbf{A}_i^T \mathbf{P}_i \left( X - \sum_{i=1}^N \mathbf{P}_i^T \mathbf{A}_i d_L \right) = \\ &= - \sum_{i=1}^N \mathbf{A}_i^T \mathbf{P}_i (X - \hat{X}). \end{aligned} \quad (\text{A.7})$$

Note that  $\mathbf{P}_i(X - \hat{X})$  is the  $i$ -th patch of the residual image, and by substituting back the definition of  $\mathbf{A}_i$ , and using the same property as before (see [9] property no. (40)) we can write the gradient as

$$\begin{aligned} \nabla_{d_L} &= - \sum_{i=1}^N (\alpha_i \otimes \mathbf{I}_{n \times n})^T \mathbf{P}_i (X - \hat{X}) = \\ &= - \sum_{i=1}^N \text{vec}(\mathbf{P}_i (X - \hat{X}) \cdot \alpha_i^T), \end{aligned} \quad (\text{A.8})$$

where  $\text{vec}(\cdot)$  denotes the vec-operator that stacks the columns of the gradient matrix into a vector, so by reshaping the above expression we get the final expression for the gradient (11)

$$\nabla_{\mathbf{D}_L} = - \sum_{i=1}^N \mathbf{P}_i (X - \hat{X}) \cdot \alpha_i^T. \quad (\text{A.9})$$

## B. Complexity Analysis

This section details the computational complexity analysis of our algorithm and compares it to previous methods.

We assume that the number of the non-zeros in every needle is limited to at most  $k$  non-zeros, and we denote by  $I$  the number of the training images. We evaluate the complexity of every outer iteration (single epoch) of our algorithm and compare it to the complexity of executing an epoch in the alternative algorithms. Every inner iteration of our algorithm (Algorithm 1 in the main paper) operates on a layer of  $N/n$  needles, while the global iteration operates on the whole dataset. The resulting computational cost of the residuals  $R_j$  for all the layers is  $O(IN(nk + n))$ ,

which is comprised of  $O(INnk)$  for computing all the  $N$  slices  $\mathbf{D}_L \alpha_i$  of all the  $I$  images, and  $O(INn)$  computations for subtracting them from the global residual. Given the residuals, every iteration applies the LARS algorithm for solving the local sparse pursuit for all the  $NI$  needles, requiring  $O(k^3 + mk^2 + nm)$  per needle [8], and  $O(IN(k^3 + mk^2 + nm) + n^2m^2)$  computations for all the  $N$  needles in all the  $I$  images. The latter term,  $n^2m^2$ , corresponds to the precomputation of the Gram matrix of the dictionary  $\mathbf{D}_L$  at every layer, which is usually negligible since it is computed once for all the needle in the layer. Next, we evaluate the complexity of reconstructing the signal. Direct computation requires  $O(IN(nk + n + k))$  operations, which is effectively  $O(INnk)$ , since this is the dominant term. The computation of the global residual requires another  $O(IN)$  operations. These last two phases are negligible compared to the sparse pursuit stage, and are therefore omitted from the final expression. Finally, the computation of the gradient is  $O(IN(n + nk))$ , and the dictionary update stage is  $O(INm)$  (updating the dictionary requires  $O(mn)$  computations and occurs  $IN/n$  times in every epoch). We summarize the above analysis in Table 1 in the main paper, and compare it to the complexity analysis of the SBDL algorithm [10]. In addition, Table 1 presents the complexity of executing an epoch of the two SGD based online algorithms that were introduced in [6], where  $q$  corresponds to the number of inner iterations of the sparse pursuit stage (performed in the Fourier domain).

The most demanding stage, in both the SBDL algorithm and in our approach, is the local sparse pursuit, which is  $O(IN(k^3 + mk^2 + nm))$ . Assuming that the needles are very sparse  $k \ll m$ , which is often the case with real-world signals, the complexity of the local sparse pursuit stage is governed by  $O(NInm)$  in both algorithms. This implies that our algorithm is of comparable order of complexity<sup>1</sup> as SBDL, which is a batch algorithm. On the other hand, the complexity of the online algorithm in [6] is dominated by the computation of the FFT in the pursuit stage, which is  $O(qImN \log(N))$ , meaning that their algorithm scales as  $O(N \log(N))$  with the global dimension of the signals, while our algorithm grows linearly.

## C. Multi-focus

In this section we address the task of multi-focus image fusion and further detail our implementation.

We denote the set of source images to fuse as  $\{Y^k\}_{k=1}^L$ , and assume that a pretrained dictionary  $\{d_i\}_{i=1}^m$  is provided. Each image  $Y^k$  is decomposed into a base component  $Y_b^k$ , which is a smooth piece-wise constant image, and an edge

<sup>1</sup>While the complexity of our algorithm is of the same order as the complexity of SBDL, we empirically demonstrated that our method converges faster (see Figure 2 in the main paper).

component  $Y_e^k$  that contains the high frequency elements:

$$Y^k = Y_b^k + Y_e^k, \quad (\text{C.1})$$

where the separation is performed by means of applying distinctive priors. The base component is usually extracted by imposing a prior which penalizes the  $l_2$  norm of its gradient. Modeling the edge component, however, is more involved and has been the subject matter of many image-processing algorithms [3, 5, 13, 14, 16, 17]. In this work, we employ the CSC model to describe the edge components, as it has shown promising results in [7].

Using the aforementioned priors, the separation of the image to its components amounts to solving the following optimization problem:

$$\min_{\Gamma_e^k, Y_b^k} \frac{1}{2} \|Y^k - \mathbf{D}_e \Gamma_e^k - Y_b^k\|_2^2 + \lambda \|\Gamma_e^k\|_1 + \mu \frac{1}{2} \|\nabla Y_b^k\|_2^2, \quad (\text{C.2})$$

where  $\Gamma_e^k$  is the sparse representation of  $Y_e^k$ , under the given convolutional dictionary  $\mathbf{D}_e$ , i.e.  $Y_e^k = \mathbf{D}_e \Gamma_e^k$ , and  $\|\nabla Y_b^k\|_2^2$  is given by

$$\|\nabla Y_b^k\|_2^2 = \|g_x * Y_b^k\|_2^2 + \|g_y * Y_b^k\|_2^2, \quad (\text{C.3})$$

where  $g_x = [-1 \ 1]$  and  $g_y = [-1 \ 1]^T$  are the horizontal and vertical gradient operators, respectively.

By taking similar steps to those presented in Section 3.1 in the main paper, we can rewrite the above optimization problem as

$$\begin{aligned} \min_{\{\alpha_j^k\}, Y_b^k} \frac{1}{2} \|Y^k - \sum_{j=1}^N \mathbf{P}_j^T \mathbf{D}_L \alpha_j^k - Y_b^k\|_2^2 \\ + \lambda \sum_{j=1}^N \|\alpha_j^k\|_1 + \mu \frac{1}{2} \|\nabla Y_b^k\|_2^2, \end{aligned} \quad (\text{C.4})$$

where  $\{\alpha_j^k\}_{j=1}^N$  are the needles which compose the sparse vector  $\Gamma_e^k$ , and  $\mathbf{D}_L$  is the local dictionary of  $\mathbf{D}_e$ .

This problem can be solved by alternating between minimizing w.r.t.  $Y_b^k$  and  $\Gamma_e^k$ , where the latter boils down to seeking for the sparse needles  $\{\alpha_j^k\}_{j=1}^N$ . To that end, the update rule of the  $\{\alpha_j^k\}_{j=1}^N$  is the set of local pursuit problems:

$$\min_{\{\alpha_j^k\}} \frac{1}{2} \|(Y^k - Y_b^k) - \sum_{j=1}^N \mathbf{P}_j^T \mathbf{D}_L \alpha_j^k\|_2^2 + \lambda \sum_{j=1}^N \|\alpha_j^k\|_1, \quad (\text{C.5})$$

which can be solved using our proposed algorithm, whereas the update rule of  $Y_b^k$  is the following least square minimization problem:

$$\min_{Y_b^k} \frac{1}{2} \|(Y^k - \mathbf{D}_e \Gamma_e^k) - Y_b^k\|_2^2 + \mu \frac{1}{2} \|\nabla Y_b^k\|_2^2. \quad (\text{C.6})$$

For solving problem C.6, we set its gradient w.r.t.  $Y_b$  to zero to obtain the following update rule:

$$Y_b^k = (I + \mu(G_x^T G_x + G_y^T G_y))^{-1} (Y^k - \mathbf{D}_e \Gamma_e^k), \quad (\text{C.7})$$

where  $G_x$  and  $G_y$  are the matrix representations of the gradient-operators.

Once these problems have been solved for all the input images  $\{Y^k\}_{k=1}^L$ , we aim to merge each set of feature maps<sup>2</sup>  $\{Z_i^k\}_{k=1}^L$  in a way that best captures the focused objects in the resulting images. For each image  $Y^k$ , we generate an activity map based on the intensity of the  $l_1$ -norm of its feature maps. More specifically, we sum pixel-wise the absolute value of its  $m$  feature maps  $\{Z_i^k\}_{i=1}^m$  to form an activity map that matches the size of the image  $N$ :

$$\tilde{\mathbf{A}}^k(u, v) = \sum_{i=1}^m \|Z_i^k(u, v)\|_1. \quad (\text{C.8})$$

To make this method more robust and less susceptible to misregistration, we convolve the above activity maps with a uniform kernel  $U_s$ , of a small support  $s \times s$ , to produce the final activity maps:

$$\mathbf{A}^k = \tilde{\mathbf{A}}^k * U_s. \quad (\text{C.9})$$

Based on the observation that a significant value in the activity map  $\mathbf{A}^k$  indicates a sharp region in the image  $Y^k$ , we then reconstruct the all-in-focus edge component by selectively assembling the most prominent regions from the feature maps based on their pixel-wise values in the corresponding activity maps:

$$Z_i^f(u, v) = Z_i^{k^*}(u, v), \quad k^* = \arg \max_k (\mathbf{A}^k(u, v)), \quad (\text{C.10})$$

where  $\{Z_i^f\}_{i=1}^m$  are the feature maps of the fused image.

Afterward, we fuse the base components, either by taking their average

$$Y_b^f = \frac{1}{N} \sum_{k=1}^L Y_b^k, \quad (\text{C.11})$$

or by nominating regions of the base components according to the maximum value in the respective activity maps, i.e.

$$Y_b^f(u, v) = Y_b^{k^*}(u, v), \quad k^* = \arg \max_k (\mathbf{A}^k(u, v)), \quad (\text{C.12})$$

where  $Y_b^f$  is the base component of the fused image. Here, we opt for the latter since it produces better results.

Finally, the fusion result  $Y^f$  is obtained by gathering its components:

$$Y^f = Y_b^f + \sum_{i=1}^m d_i * Z_i^f. \quad (\text{C.13})$$

<sup>2</sup>This set of feature maps refers to the  $i$ -th feature maps of the input images.



## D. Experiments

In this section we turn to describe the implantation details of our experiments.

Throughout all our experiments we used a local dictionary composed of  $m = 81$  filters, each of size  $8 \times 8$ . In addition, we used the LARS algorithm [1] for solving the local sparse pursuit stage (problem (8) in the main paper).

**Run time comparison of the batch methods (Figure 2 in the main paper).** We used  $\lambda = 1$  and the Fruit dataset [2] for training the dictionaries. The dataset contains 10 images of size  $100 \times 100$  pixels. As a preprocessing step, we mean-subtracted the images. The mean was computed by convolving each image with an  $8 \times 8$  uniform kernel and subtracting the result from the original image. Additionally, we used the ADAM algorithm in the initial 30 iterations, with  $\eta = 0.02$ , and set the ADAM parameters in accordance with the authors' recommendation:  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ , and  $\epsilon = 10^8$ . Subsequent iterations applied the Momentum algorithm with  $\eta = 10^{-7}$  and  $\gamma = 0.8$  until convergence<sup>3</sup>.

**Run time comparison of the online algorithms (Figure 3 in the main paper).** For training the dictionary we used the MIRFLICKR-1M dataset [4], where we sampled (without recurrence) 40 images for training and an additional 5 images for testing. The images were cropped to reduce their size from  $512 \times 512$  to  $256 \times 256$  pixels in both the training and testing sets to expedite the computation. In addition, we divided the images by 255 and mean-subtracted them as previously described. In this experiment, we used  $\lambda = 0.1$  for all methods. For our method we set the learning rate to  $\eta = 0.1$ , with learning rate decay<sup>4</sup> of  $\tilde{\eta} = \eta / (1 + 3/t)$  updated every 5 epochs. Lastly, we setup the momentum parameter  $\gamma = 0.8$ .

**Multi-focus image fusion (Section 7.3 in the main paper).** The dictionaries of both methods were pretrained on the Fruit dataset [2]. In addition, we set the coefficients described in Equation (C.2) to  $\lambda = 1$  and  $\mu = 5$  for the sparse pursuit (C.5) and the base-image extraction (C.6) stages, respectively, and alternate between the stages at each iteration. In practice, convergence was achieved within 2-4 iterations of alternating between these two stages. For reconstructing the image Barbara we used a reconstruction kernel  $U_s$  (Equation (C.9)) of size  $9 \times 9$ , and for the image Butterfly we used a reconstruction kernel  $U_s$  of size  $8 \times 8$ .

For reconstructing the colored image Bird in the Lab color space, we built the activity maps  $A_k$  based on their

<sup>3</sup>All our notations are in accordance with those presented in [12]

<sup>4</sup>t here denotes the numbers of the epochs.

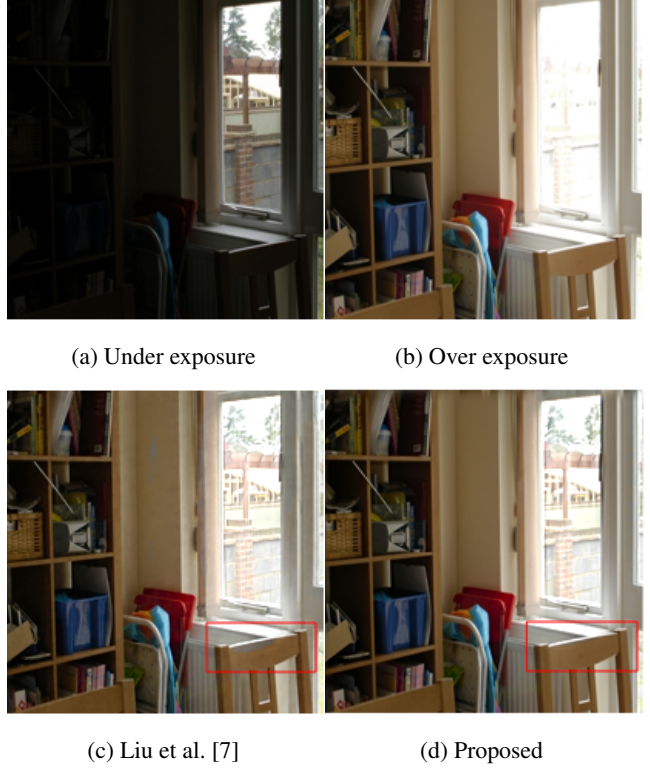


Figure 1: Multi-exposure fusion comparison of real images. The red boxes indicate artifact suppression characteristics of the proposed method.

L channel only, with kernel  $U_s$  of size of  $14 \times 14$ . Then, we reconstructed each channel from the Lab color space by selecting regions based on the maximum pixel-wise value of the activity maps. Finally for fusing the image-pair Clocks we used a kernel  $U_s$  of size  $9 \times 9$ .

## E. Additional Applications

We turn to demonstrate the benefits of LoBCoD on two additional applications: Multi-exposure image fusion and salt-and-pepper noise removal from text images.

**Multi-exposure image fusion.** To showcase the versatility of the multi-focus LoBCoD application, we adapt it to the task of multi-exposure image-fusion. The fusion is preformed in the same manner as in the case of multi-focus fusion (Section C), with a slight modification of handling the base components. Here, we fuse the base components by taking their weighted sum:

$$Y_b^f = \sum_{k=1}^L a^k Y_b^k. \quad (\text{E.1})$$

Method	PSNR [dB]
Noisy test set	13.37
Wohlberg [15] trained on clean training set	19.83
Plaut et al. [11] trained on clean training set	21.60
Proposed trained on clean training set	20.40
Wohlberg [15] trained on noisy test set	19.62
Plaut et al. [11] trained on noisy test set	22.30
Proposed trained on noisy test set	23.64

Table 1: Average PSNR comparison between our method and the algorithms presented in [11] and [15]. The noisy images are corrupted by converting 10% of the image-pixels.

All other implementation details, including the parameters  $\lambda$  and  $\mu$  remain the same as in the multi-focus case. For fusing the image Window (Figure 1) we used weights of 0.9 and 0.1 for the over and under exposed images, respectively, as it yields a brighter combined image. For  $U_s$  we chose a reconstruction kernel of size of  $6 \times 6$ .

Figure 1 compares our results with the results of an adapted version of Liu et al. [7], demonstrating LoBCoD’s superior performance.

**Salt-and-pepper text image denoising.** Salt-and-pepper noise, also known as impulse noise, is characterized by sparsely occurring white or black (maximum or minimum value) pixels in an image. Here, unlike the inpainting task, we do not assume prior knowledge of the noise mask and we rely solely on the sparse prior to reconstruct the corrupted image. Thus, for a given test image, we aim to solve the following minimization problem:

$$\min_{Z_i} \sum_{i=0}^m \|Z_i\|_1 \quad s.t. \quad \|X - \sum_{i=0}^m d_i * Z_i\|_2^2 < \epsilon. \quad (\text{E.2})$$

In problem E.2 we adopt the approach presented in [11, 15] and augment the dictionary  $\{d_i\}_{i=1}^m$  with a single impulse filter  $d_0 = [1 \ 0 \ 0 \ \dots \ 0]^T$  to represent the impulse noise. For solving problem (E.2), we use the Lagrangian formulation as described in the main paper but with a small modification of assigning a different  $\lambda$  for image reconstruction and noise representation. This is because the sparsity ratio of an image presented by the dictionary  $\{d_i\}_{i=1}^m$  depends on the complexity of the image, while the sparsity ratio of the noise depends on the noise level (the amount of corrupted pixels). Thus, the objective for our minimization can be

written as follows:

$$\min_{Z_i} \frac{1}{2} \|X - \sum_{i=1}^m d_i * Z_i - d_0 * Z_0\|_2^2 + \lambda_1 \sum_{i=1}^m \|Z_i\|_1 + \lambda_0 \|Z_0\|_1. \quad (\text{E.3})$$

We minimize the objective (E.3) by splitting it into two minimization sub-problems, one minimizes over  $Z_0$  leaving  $\{Z_i\}_{i=1}^m$  fixed, and the other minimizes over  $\{Z_i\}_{i=1}^m$  with  $Z_0$  fixed. Then we alternate between these two minimization sub-problems for several iterations. In practice, 2-4 iterations were sufficient for convergence. For reconstructing the restored image, we use only the dictionary  $\{d_i\}_{i=1}^m$  with its coefficients and zero the coefficient of the impulse atom.

We evaluate our proposed algorithm using text images taken from the dataset in [11], which is a subset of scanned pages taken from the book Aristotle’s Nicomachean Ethics. Each training and test set are composed of 16 pages from the book, each of size  $493 \times 383$  pixels. The images are gray scale images, normalized to the range  $[0, 1]$ . Since the the dataset assigns a value of zero to text characters and one to the image background, we employ a pre-processing step of inverting the images as described in [11], and uninvert the results to obtain black-over-white images at the end.

The images in the test set were corrupted by randomly inverting 10% of the pixels in every image, setting 5% of the pixels to black and 5% of the pixels to white. Figure 2(a) shows a portion of an original image taken from the test set and Figure 2(b) shows its corrupted version. We apply our algorithm to reconstruct the test set both in the case of training the dictionary on the clean training set and training on the noisy test images, and compare our results to the methods presented in [11] and [15] using their published code. In all evaluated methods, the dictionary was comprised of 100 atoms, each atom of size  $11 \times 11$  pixels. Additionally, when training each method on the noisy test set, we include a pruning procedure of noisy atoms, as presented in [11].

Table 1 presents the resulting average PSNR of the different algorithms on the test set. The comparison illustrates that while [11] performs better when training on clean images, the proposed method achieves better results when training on the corrupted test images. Figure 2 compares the various algorithms when applied to a single test example. Here, the advantage of our approach is reiterated in the case of training using noisy test images.



Figure 2: Denoising performance comparison on a text image. (a) the original image; (b) the noisy image. PSNR: 13.53; (c) denoising method [15], dictionary trained on clean images. PSNR: 19.23; (d) denoising method [11], dictionary trained on clean images. PSNR: 20.68; (e) denoising using the proposed algorithm, dictionary trained on clean images. PSNR: 20.25; (f) denoising method [15], dictionary trained on noisy images. PSNR: 18.51; (g) denoising method [11], dictionary trained on noisy images. PSNR: 19.91; (h) denoising using the proposed algorithm, dictionary trained on noisy images. PSNR: 22.81.

## References

- [1] B. Efron, T. Hastie, I. Johnstone, R. Tibshirani, et al. Least angle regression. *The Annals of statistics*, 32(2):407–499, 2004.
- [2] R. Fergus, M. D. Zeiler, G. W. Taylor, and D. Krishnan. Deconvolutional networks. In *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 00, pages 2528–2535, 2010.
- [3] R. Gao and S. A. Vorobyov. Multi-focus image fusion via coupled sparse representation and dictionary learning. *arXiv preprint arXiv:1705.10574*, 2017.
- [4] M. J. Huiskes, B. Thomee, and M. S. Lew. New trends and ideas in visual concept detection: the mir flickr retrieval evaluation initiative. In *Proceedings of the international conference on Multimedia information retrieval*, pages 527–536. ACM, 2010.
- [5] H. Li, B. Manjunath, and S. K. Mitra. Multisensor image fusion using the wavelet transform. *Graphical models and image processing*, 57(3):235–245, 1995.
- [6] J. Liu, C. Garcia-Cardona, B. Wohlberg, and W. Yin. First-and second-order methods for online convolutional dictionary learning. *SIAM Journal on Imaging Sciences*, 11(2):1589–1628, 2018.
- [7] Y. Liu, X. Chen, R. K. Ward, and Z. J. Wang. Image fusion with convolutional sparse representation. *IEEE signal processing letters*, 23(12):1882–1886, 2016.
- [8] J. Mairal, F. Bach, J. Ponce, et al. Sparse modeling for image and vision processing. *Foundations and Trends® in Computer Graphics and Vision*, 8(2-3):85–283, 2014.
- [9] T. P. Minka. Old and new matrix algebra useful for statistics. See [www.stat.cmu.edu/minka/papers/matrix.html](http://www.stat.cmu.edu/minka/papers/matrix.html), 2000.
- [10] V. Pappas, Y. Romano, J. Sulam, and M. Elad. Convolutional dictionary learning via local processing. In *ICCV*, pages 5306–5314, 2017.
- [11] E. Plaut and R. Giryes. A greedy approach to 0,-based convolutional sparse coding. *SIAM Journal on Imaging Sciences*, 12(1):186–210, 2019.
- [12] S. Ruder. An overview of gradient descent optimization algorithms. *arXiv preprint arXiv:1609.04747*, 2016.
- [13] S. Savić. Multifocus image fusion based on empirical mode decomposition. In *Twentieth International Electro technical and Computer Science Conference*, 2011.
- [14] W. Wang and F. Chang. A multi-focus image fusion method based on laplacian pyramid. *JCP*, 6(12):2559–2566, 2011.
- [15] B. Wohlberg. Convolutional sparse representations as an image model for impulse noise restoration. In *2016 IEEE 12th Image, Video, and Multidimensional Signal Processing Workshop (IVMSP)*, pages 1–5. IEEE, 2016.

- [16] B. Yang and S. Li. Multifocus image fusion and restoration with sparse representation. *IEEE Transactions on Instrumentation and Measurement*, 59(4):884–892, 2010.
- [17] B. Yang and S. Li. Pixel-level image fusion with simultaneous orthogonal matching pursuit. *Information fusion*, 13(1):10–19, 2012.