

000
001
002
003
004
005
006
007
008
009
010
011

Abstract

Convolutional sparse coding (CSC) is a promising direction for unsupervised learning in computer vision. In contrast to recent supervised methods, CSC allows for convolutional image representations to be learned that are equally useful for high-level vision tasks and low-level image reconstruction and can be applied to a wide range of tasks without problem-specific retraining. Due to their extreme memory requirements, however, existing CSC solvers have so far been limited to low-dimensional problems and datasets using a handful of low-resolution example images at a time.

In this paper, we propose a new approach to solving CSC as a consensus optimization problem, which lifts these limitations. By learning CSC features from large-scale image datasets for the first time, we achieve significant quality improvements in a number of imaging tasks. Moreover, the proposed method enables new applications in high-dimensional feature learning that has been intractable using existing CSC methods. This is demonstrated for a variety of reconstruction problems across diverse problem domains, including 3D multispectral demosaicing and 4D light field view synthesis.

012
013

1. Introduction

Natural image statistics lie at the core of a wide variety of discriminative and generative computer vision tasks. In particular, convolutional image representations have proven essential for supervised learning using deep neural networks – the de-facto state-of-the-art for many high-level vision tasks [20, 28, 27, 13]. While these models are successful for supervised discriminative problems, the same architectures do not easily transfer to generative tasks.

Generative models have some significant advantages over discriminative models for low level vision and image reconstruction tasks. The most important distinction is that generative approaches learn models of the data that can act as priors for a wide range of reconstruction tasks without retraining, while discriminative methods learn specific reconstruction tasks, and cannot be easily applied to other

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100

tasks. As a consequence patch-based sparse coding techniques [7, 22, 1] have been very popular for low-level tasks such as denoising, inpainting, demosaicing, deconvolution and similar problems [11, 33, 29, 23, 21, 2]. Unfortunately, patch-based dictionaries are highly redundant because they have to capture all shifted copies of the sparsifying filters.

Introduced as a model for receptive fields in human vision [25], convolution sparse coding (CSC) [14, 17, 31, 32] has been demonstrated to remove much of the overhead of patch-based sparse coding by using a convolution image formation model for a range of different applications [11, 33, 29, 23, 21, 2]. CSC techniques are fast, because many implementations efficiently perform convolutions in the frequency domain [5, 6, 15]. While fast, existing CSC approaches are not scalable due to their extreme memory requirements (Fig. 3). For example, existing methods would require terabytes of physical memory for learning light field data with only 100 examples (Sec. 4), and datasets comparable to ImageNet would require petabytes of memory. As a result, it has been intractable to learn convolutional filters from large datasets, and to apply CSC to high-dimensional image reconstruction problems that arise in 3D video, 3D

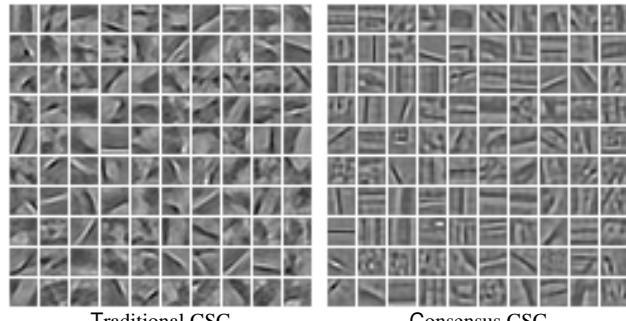
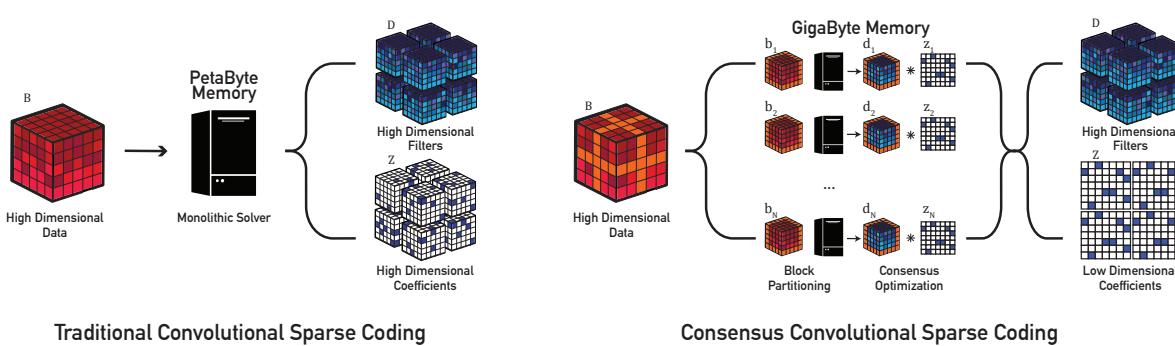


Figure 1: Large-scale unsupervised feature learning. Left: Convolutional features from [15] can only be learned from a handful of example images since existing CSC methods are limited by memory. Right: CCSC overcomes these limitations, and allows to learn features on ImageNet [9]. These features contain less specialized structures, leading to significant improvements across a variety of vision tasks.



Traditional Convolutional Sparse Coding

Consensus Convolutional Sparse Coding

Figure 2: Illustration of traditional CSC (left) and the proposed CCSC (right). CCSC lifts the prohibitive memory limitations of existing algorithms by breaking large, high dimensional datasets into tractable subproblems which can be efficiently solved with low memory footprint.

multispectral, or 4D light field image processing.

In this paper, we revisit *unsupervised, generative* learning using CSC, and propose a consensus-based optimization framework that makes CSC tractable on large-scale datasets, and enables high-dimensional feature learning. We call our approach consensus convolutional sparse coding (CCSC). CCSC splits a single large-scale problem into a set of smaller sub-problems that fit into available memory resources. Due to the convex nature of the problem and the enforced consensus between the sub-problems, global convergence is guaranteed. We demonstrate convolutional dictionary learning on datasets that are orders of magnitude larger than what has previously been possible, and show that the resulting sparsifying filters are, in fact, different from those learned from smaller datasets (Fig. 1). Moreover, we show that these new features also lead to significant improvements in a variety of image reconstruction tasks. To validate the proposed method for high-dimensional data, we evaluate CCSC on a number of high-dimensional reconstruction problems that are intractable for existing CSC solvers. In particular, we make the following contributions:

- We derive a consensus optimization method that enables convolutional sparse coding problems of arbitrary size with limited memory to be solved efficiently.
- We extend traditional CSC to allow for non-convolutional data dimensions, greatly reducing memory requirements for high-dimensional datasets.
- We verify the scalability of CCSC by learning from large-scale 2D datasets as well as from several high-dimensional datasets.
- We show that the features learned on large-scale datasets are more general, and lead to better reconstructions than existing methods.
- We evaluate CCSC using several high-dimensional reconstruction problems across diverse problem do-

mains, including 3D multispectral demosaicing, 3D video deblurring, and 4D light field view synthesis.

Finally, the full source code will be made available online for evaluation and improvements in the future.

2. Mathematical Framework

Traditionally, convolutional sparse coding is formulated as the following optimization problem

$$\underset{\mathbf{d}, \mathbf{z}}{\operatorname{argmin}} \sum_{j=1}^J \frac{1}{2} \|\mathbf{b}^j - \sum_{w=1}^W \mathbf{d}_w * \mathbf{z}_w^j\|_2^2 + \beta \sum_{w=1}^W \|\mathbf{z}_w^j\|_1 \quad (1)$$

subject to $\|\mathbf{d}_w\|_2^2 \leq 1 \quad \forall w \in \{1, \dots, W\}$,

where each example image \mathbf{b}^j is represented as the sum of sparse coefficient feature maps \mathbf{z}_w^j convolved with filters \mathbf{d}_w of fixed spatial support. The superscripts indicate the example index $j = 1 \dots J$, and the subscripts indicate the filter/coefficient map index $w = 1 \dots W$. The variables $\mathbf{b}^j \in \mathbb{R}^D$ and $\mathbf{z}_w^j \in \mathbb{R}^D$ are vectorized images and feature maps, respectively, $\mathbf{d}_w \in \mathbb{R}^M$ represents the vectorized m -dimensional filters, and $*$ is the m -dimensional convolution operating on the vectorized inputs. The constraint on \mathbf{d}_w ensures the dictionary does not absorb all of the system's energy.

To solve Eq. (1) we first reformulate it as an unconstrained optimization problem, following [15]. Absorbing the constraint in an additional indicator penalty $\text{ind}_C(\cdot)$ for each filter, defined on the convex set of constraints $C = \{x \mid \|Sx\|_2^2 \leq 1\}$, where S projects the filter onto its spatial support, yields

$$\underset{\mathbf{d}, \mathbf{z}}{\operatorname{argmin}} \frac{1}{2} \sum_j (\|\mathbf{b}^j - \mathbf{Z}^j \mathbf{d}\|_2^2 + \beta \|\mathbf{Z}^j\|_1 + \text{ind}_C(\mathbf{d})). \quad (2)$$

Here, $\mathbf{d} = [\mathbf{d}_1^T \dots \mathbf{d}_W^T]^T$ and $\mathbf{Z}^j = [\mathbf{Z}_1^j \dots \mathbf{Z}_W^j]$ is a concatenation of Toeplitz matrices, each one expressing the con-

216 convolution with sparse coefficient map \mathbf{z}_w^j . Note that we can
 217 express the convolutional term from Eq. (1) in this way because
 218 convolution is a commutative operator. Eliminating
 219 the sum over the examples (index J) by stacking the vec-
 220 torized images in $\mathbf{b}' = [\mathbf{b}_1^T \dots \mathbf{b}_J^T]^T$ and coefficient maps
 221 $\mathbf{Z}' = [\mathbf{Z}^1 \dots \mathbf{Z}^J]^T$ accordingly results in
 222

$$\operatorname{argmin}_{\mathbf{d}, \mathbf{z}} \frac{1}{2} \|\mathbf{b}' - \mathbf{Z}' \mathbf{d}\|_2^2 + \beta \|\mathbf{Z}'\|_1 + \operatorname{ind}_C(\mathbf{d}). \quad (3)$$

223 We jointly solve for both the filters \mathbf{d} and coefficient maps
 224 \mathbf{z} in Equation 3 using a coordinate descent approach [15]
 225 that alternates between updates to \mathbf{d} and \mathbf{z} while keeping
 226 the other fixed (described later in Alg. 2). For this spatial
 227 formulation, the filters can be represented in a memory-
 228 efficient way, due to their small spatial support. However,
 229 the full set of coefficients \mathbf{z}_w must be stored which incurs
 230 an enormous memory footprint. Furthermore, convolutions
 231 in the spatial domain are computationally expensive.
 232

233 Recent work [5, 6, 18, 15] has demonstrated that Eq. (3)
 234 can be solved efficiently in the frequency domain by ap-
 235 plying Parseval’s theorem, which states that the energy
 236 of a signal is equivalent to that of its Fourier transform
 237 up to a constant. In this frequency domain formulation,
 238 the previously costly spatial convolutions become efficient
 239 Hadamard (component-wise) products. Although computa-
 240 tionally efficient, the Fourier formulation still requires fre-
 241 quency representations over the full domain of all frequen-
 242 cies to be held in memory, both for filters and coefficient
 243 maps. The size of the coefficient maps grows linearly with
 244 the number of filters and images, but exponentially with the
 245 dimensionality. For these reasons, classical convolutional
 246 sparse coding, and especially its efficient Fourier formu-
 247 lation, do not scale beyond 2D images and small training
 248 datasets.
 249

250 In the following, we derive a consensus optimization
 251 method for CSC, allowing to split large-scale and high-
 252 dimensional CSC into smaller sub-problems, each of which
 253 can be solved with a limited memory budget. Furthermore,
 254 the individual sub-problems can be solved efficiently using
 255 the Fourier-domain formulation, and in a distributed fash-
 256 ion using parallel workers. Consensus optimization makes
 257 CSC tractable for large problems sizes, which we verify by
 258 learning from large-scale and high-dimensional datasets.
 259

260 2.1. Consensus Optimization

261 To account for large, high-dimensional datasets, we split
 262 the problem of learning from the entire dataset \mathbf{b}' into learn-
 263 ing from smaller subsets which can be solved individu-
 264 ally with modest memory and computational requirements.
 265 Specifically, we partition the data vector \mathbf{b}' and their cor-
 266 responding sparse feature matrix \mathbf{Z}' across all of the exam-

267 ples¹ into N blocks arranged by rows,
 268

$$\mathbf{b}' = \begin{bmatrix} \mathbf{b}_1 \\ \vdots \\ \mathbf{b}_N \end{bmatrix}, \quad \mathbf{Z}' = \begin{bmatrix} \mathbf{Z}_1 \\ \vdots \\ \mathbf{Z}_N \end{bmatrix}, \quad (4)$$

269 with $\mathbf{b}_i \in \mathbb{R}^{B_i}$ and $\mathbf{Z}_i \in \mathbb{R}^{B_i \times MW}$, where $\sum_{i=1}^N B_i =$
 270 JD . Here, \mathbf{b}_i represents the i^{th} data block along with its
 271 respective filters \mathbf{Z}_i . In the following we first demonstrate
 272 how to solve Eq. (3) using this block splitting with respect
 273 to the filters \mathbf{d} , and subsequently for the coefficients \mathbf{z} .
 274

275 2.1.1 Filter Subproblem

276 Using the partition from Eq. (4), we can solve Eq. (3) for \mathbf{d}
 277 for a given \mathbf{Z}' as follows
 278

$$\begin{aligned} & \operatorname{argmin}_{\mathbf{d}} \frac{1}{2} \sum_{i=1}^N \|\mathbf{b}_i - \mathbf{Z}_i \mathbf{d}\|_2^2 + \operatorname{ind}_C(\mathbf{d}) \\ & \Leftrightarrow \operatorname{argmin}_{\mathbf{y}} \frac{1}{2} \sum_{i=1}^N \|\mathbf{b}_i - \mathbf{Z}_i \mathbf{d}_i\|_2^2 + \operatorname{ind}_C(\mathbf{y}) \\ & \text{subject to } \mathbf{d}_i - \mathbf{y} = 0 \quad \forall i \in \{1, \dots, N\}. \end{aligned} \quad (5)$$

279 This is a convex problem in the global consensus form [3].
 280 Introducing local variables \mathbf{d}_i allows us to turn the joint ob-
 281 jective from the first row of Eq. (5), which cannot be split
 282 due to the joint variable \mathbf{d} , into separable terms that can be
 283 split during the optimization. This also facilitates the han-
 284 dling of the i -th set $(\mathbf{b}_i, \mathbf{Z}_i, \mathbf{d}_i)$ independently by parallel
 285 workers. The shared global variable $\mathbf{y} \in \mathbb{R}^{MW}$ introduced
 286 as a slack variable enables solving Eq. (5) using the Alter-
 287 nate Direction Method of Multipliers (ADMM) [3], which
 288 we derived from the augmented Lagrangian
 289

$$\begin{aligned} \mathcal{L}(\mathbf{d}_1 \dots \mathbf{d}_N, \mathbf{y}, \lambda_1 \dots \lambda_N) = & \sum_{i=1}^N \frac{1}{2} \|\mathbf{b}_i - \mathbf{Z}_i \mathbf{d}_i\|_2^2 \\ & + \operatorname{ind}_C(\mathbf{y}) + \lambda_i^T (\mathbf{d}_i - \mathbf{y}) + \frac{\rho}{2} \|\mathbf{d}_i - \mathbf{y}\|_2^2, \end{aligned} \quad (6)$$

290 where λ_i is a set of a Lagrange multipliers for each of the
 291 N consensus constraints. ADMM alternately minimizes
 292 Eq. (6) with respect to all of its variables, yielding Alg. 1.
 293

294 Line 5 uses the average $\bar{\mathbf{d}}^{k+1} = \frac{1}{N} \sum_{i=1}^N \mathbf{d}_i^{k+1}$ and
 295 $\bar{\lambda}^k = \frac{1}{N} \sum_{i=1}^N \lambda_i^k$ as a notational shortcut. It becomes clear
 296 that the subproblems in the first inner for-loop around Line 3
 297 are now independent of each other. The N subproblems can
 298 be solved on a single machine sequentially, or in parallel on
 299 up to N workers, each worker i handling only the i -th block
 300 of data. After the parallel solve a global synchronization
 301 step in Line 5 fuses all individual filter dictionaries, while
 302

303 ¹Please see the supplemental for other splitting strategies.
 304

324
 325 **Algorithm 1** ADMM for the Filters \mathbf{d}
 326 1: **while** Not Converged **do**
 327 2: **for** $i = 1$ to N **do**
 328 3: $\mathbf{d}_i^{k+1} = \operatorname{argmin}_{\mathbf{d}_i} \frac{1}{2} \|\mathbf{b}_i - \mathbf{Z}_i \mathbf{d}_i\|_2^2 + \frac{\rho}{2} \|\mathbf{d}_i - \mathbf{y}^k + \lambda_i^k\|_2^2$
 329 4: **end for**
 330 5: $\mathbf{y}^{k+1} = \operatorname{argmin}_{\mathbf{y}} \operatorname{ind}_C(\mathbf{y}) + \frac{N\rho}{2} \|\mathbf{y} - \bar{\mathbf{d}}^{k+1} - \bar{\lambda}^k\|_2^2$
 331 6: **for** $i = 1$ to N **do**
 332 7: $\lambda_i^{k+1} = \lambda_i^k + \mathbf{d}_i^{k+1} - \mathbf{y}^{k+1}$
 333 8: **end for**
 334 9: **end while**
 335 10: $\mathbf{d} = \mathbf{y}^{k+1}$

338
 339 enforcing the constraint $C = \{x \mid \|Sx\|_2^2 \leq 1\}$. Line 7 updates the Lagrange multipliers for each data-block based on
 340 the running error of the fused filters. In the following, we
 341 define the individual subproblems of Alg. 1 in detail.
 342

343 Line 3 is a least-squares problem with the solution

$$\mathbf{d}_i^{k+1} = (\mathbf{Z}_i^\dagger \mathbf{Z}_i + \rho \mathbb{I})^{-1} (\mathbf{Z}_i^\dagger \mathbf{b}_i + \rho (\mathbf{y}^k - \lambda_i^k)), \quad (7)$$

344
 345 where \cdot^\dagger denotes the conjugate transpose, and \mathbb{I} denotes
 346 the identity matrix. As described in [5, 6, 15] one can
 347 find a variable reordering which makes $(\mathbf{Z}_i^\dagger \mathbf{Z}_i + \rho \mathbb{I})$ block-
 348 diagonal which we directly invert using Cholesky factorization
 349 for the individual blocks, in parallel. The update in
 350 Line 5 of Alg. 1 is in the form of a proximal operator for
 351 which a rich body of literature exists [26]. Specifically, it is
 352 $\mathbf{y}^{k+1} = \operatorname{prox}_{\frac{1}{N\rho}}(\bar{\mathbf{d}}^{k+1} + \bar{\lambda}^k)$, with
 353

$$\operatorname{prox}_{\theta \operatorname{ind}_C(\cdot)}(\mathbf{v}) = \begin{cases} \frac{\mathbf{S}\mathbf{v}}{\|\mathbf{S}\mathbf{v}\|_2} & : \|\mathbf{S}\mathbf{v}\|_2^2 \geq 1 \\ \mathbf{S}\mathbf{v} & : \text{else} \end{cases} \quad \text{Projection} \quad (8)$$

2.1.2 Coefficient Subproblem

361 The coefficient subproblem can be written as

$$\begin{aligned} & \operatorname{argmin}_{\mathbf{z}} \frac{1}{2} \|\mathbf{b}' - \mathbf{D}' \mathbf{z}\|_2^2 + \beta \|\mathbf{z}\|_1 \\ & \Leftrightarrow \operatorname{argmin}_{\mathbf{z}_i} \frac{1}{2} \sum_{i=1}^N \|\mathbf{b}_i - \mathbf{D} \mathbf{z}_i\|_2^2 + \beta \|\mathbf{z}_i\|_1. \end{aligned} \quad (9)$$

363
 364 The sparse coefficient maps \mathbf{z} can be solved analogous to
 365 the filters \mathbf{d} . This is a result of the convolution from Eq. (1)
 366 being commutative, which allows to rewrite $\mathbf{Z}' \mathbf{d} = \mathbf{D}' \mathbf{z}$
 367 in Eq. (3), with \mathbf{D}' is a block diagonal matrix with $\mathbf{D} =$
 368 $[\mathbf{D}_1 \dots \mathbf{D}_W]$ repeated along its diagonal J times, and $\mathbf{z} =$
 369 $[\mathbf{z}^1 \dots \mathbf{z}^J]^T$ and $\mathbf{z}^j = [\mathbf{z}_1^j \dots \mathbf{z}_W^j]^T$. Hence, when solving
 370 for \mathbf{z} , we can follow the recipe from the previous section,
 371 using the same block partition. The resulting algorithm can
 372 be found in the supplemental material.

2.1.3 Joint Optimization

378
 379 The previous paragraphs describe optimization methods
 380 for solving the joint objective from Eq. (1) for \mathbf{d} and \mathbf{z} .
 381 We solve for both unknowns jointly by solving the bi-
 382 convex optimization problem using coordinate descent, fol-
 383 lowing [5, 15].

Algorithm 2 Large Scale CSC Learning

386
 387 1: Initialize parameters $\rho_{\mathbf{d}} \in \mathbb{R}^+$, $\rho_{\mathbf{z}} \in \mathbb{R}^+$
 388 2: Initialize variables $\mathbf{d}^0, \mathbf{z}^0, \lambda_{\mathbf{d}}^0, \lambda_{\mathbf{z}}^0, \beta$.
 389 3: **repeat**{Outer Iterations}
 390 4: **Filter Update:**
 391 $\mathbf{d}^k, \lambda_{\mathbf{d}}^k \leftarrow$ Solve with Alg. 1 and $\rho = \rho_{\mathbf{d}}$, $\lambda = \lambda_{\mathbf{d}}^{k-1}$
 392 5: **Coefficient Update:**
 393 $\mathbf{z}^k, \lambda_{\mathbf{z}}^k \leftarrow$ Detailed in supplemental $\rho = \rho_{\mathbf{z}}$, $\lambda = \lambda_{\mathbf{z}}^{k-1}$
 394 6: **until** No more progress in both directions.

397
 398 The respective Lagrange multipliers are initialized with
 399 those from the previous iteration. ρ is a parameter of the
 400 Lagrangian which intuitively is the step size enforcing the
 401 Lagrangian step. For any positive ρ , the primal residual
 402 ($\mathbf{d}_i - \mathbf{y}$) converges to zero, yielding an optimal solution.
 403 We refer to [3] for a detailed discussion and proof of convergence.
 404 Running the sub-step algorithms for a fixed number
 405 of P steps achieved good progress in the coordinate descent
 406 step. We terminate the execution when neither sub-step can
 407 further decrease the objective.

2.2 Non-Convolutional Dimensions

408
 409 Above, we have considered all dimensions of the ex-
 410 ample data \mathbf{b} to be convolutional. However, some image
 411 modalities exist only at very low resolution, e.g. the color
 412 dimension of an RGB image. In these cases it is common
 413 that no convolutional structure can be found. We represent
 414 non-convolutional dimensions by introducing an additional
 415 replication operator $\operatorname{Rep}(\cdot)$ which repeats the sparse coef-
 416 ficient maps, that do not contain the non-convolutional di-
 417 mensions, along the missing dimensions. The original con-
 418 volutional sparse coding problem from Eq. 1 becomes
 419

$$\begin{aligned} & \operatorname{argmin}_{\mathbf{d}, \mathbf{z}} \sum_{j=1}^J \frac{1}{2} \|\mathbf{b}^j - \sum_{w=1}^W \mathbf{d}_w * \operatorname{Rep}(\mathbf{z}_w^j)\|_2^2 + \beta \sum_{w=1}^W \|\mathbf{z}_w^j\|_1 \\ & \text{subject to } \|\mathbf{d}_w\|_2^2 \leq 1 \quad \forall w \in \{1, \dots, W\}, \end{aligned} \quad (10)$$

420
 421 For example, considering a single dimension with length
 422 $\mu = 3$ for RGB image data, $\operatorname{Rep}(\cdot)$ expands the 2D feature-
 423 maps to the full three-channel data by replicating the fea-
 424 ture map 3 times along the 3rd dimension. The convolution
 425 operator is still a 2D convolution, but with full color RGB
 426 filters. In Eq. (3), the operator $\operatorname{Rep}(\cdot)$ can be represented
 427
 428
 429
 430
 431

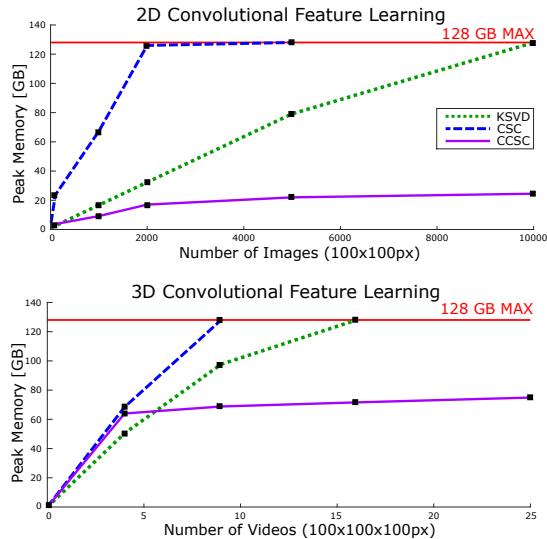


Figure 3: Memory Consumption for large 2D image datasets (top) and video data (bottom). CSC (blue) as well as popular patch-based coding methods (green) become infeasible with increasing size of the dataset (top plot). This effect is even more significant in higher dimensions (bottom plot). Note the very small number of example videos in the bottom plot. This paper proposed CCSC (magenta) which makes convolutional sparse coding suitable learning from large-scale datasets as well as higher dimensional data.

by an additional matrix $\mathbf{P} = [\mathbb{I}_1 \dots \mathbb{I}_\mu]^T$ such that \mathbf{D} and \mathbf{PZ} are then of complimentary dimensions. Redefining the coefficient matrix as $\tilde{\mathbf{Z}} = \mathbf{PZ}$, the described Alg. 1 and 2 generalize to this setting. \mathbf{P} being stacked identity matrices, the efficient inverse from Eq. (7) can be applied.

3. Memory and Complexity Analysis

This section analyzes the memory and runtime of the proposed approach. The consensus optimization from the previous section enables splitting CSC problems of arbitrary size into subproblems that fit into physical memory. Fig. 3 shows the memory consumption of the proposed CCSC approach compared to existing CSC [15], as well as classic patch-based sparse coding [1]. Even on a machine with 128 GB of physical memory these existing methods become infeasible for learning from medium datasets in 2D, and fail for small data-sets in higher-dimensions. CCSC makes large-scale convolutional sparse coding feasible by efficiently solving for smaller subproblems with memory requirements that scale slowly with increase in dataset size and dimensions. However, splitting the CSC problem comes at the cost of increased iterations which are necessary to enforce consensus between local variables.

Each subproblem can now be solved sequentially or in

Method	Cost (in flops)		
	$PJ \cdot (\underbrace{WD}_{\text{Conjugate gradient}} + \underbrace{WDM}_{\text{Spatial convolutions}} + \underbrace{WD}_{\text{Shrinkage}}$	$PJ \cdot (\underbrace{W^3D}_{\text{Linear systems}} + \underbrace{WD \log(D)}_{\text{FFTs}} + \underbrace{WD}_{\text{Shrinkage}}$	$\underbrace{W^3D + (P-1)W^2D}_{\text{Linear systems}} + PJ \cdot (\underbrace{WD \log(D)}_{\text{FFTs}} + \underbrace{WD}_{\text{Shrinkage}})$
Zeiler et al. [31]	$PJ \cdot (\underbrace{WD}_{\text{Conjugate gradient}} + \underbrace{WDM}_{\text{Spatial convolutions}} + \underbrace{WD}_{\text{Shrinkage}}$	$PJ \cdot (\underbrace{W^3D}_{\text{Linear systems}} + \underbrace{WD \log(D)}_{\text{FFTs}} + \underbrace{WD}_{\text{Shrinkage}}$	$\underbrace{W^3D + (P-1)W^2D}_{\text{Linear systems}} + PJ \cdot (\underbrace{WD \log(D)}_{\text{FFTs}} + \underbrace{WD}_{\text{Shrinkage}})$
Bristow et al. [5, 6]	$PJ \cdot (\underbrace{W^3D}_{\text{Linear systems}} + \underbrace{WD \log(D)}_{\text{FFTs}} + \underbrace{WD}_{\text{Shrinkage}}$	$PJ \cdot (\underbrace{W^3D}_{\text{Linear systems}} + \underbrace{WD \log(D)}_{\text{FFTs}} + \underbrace{WD}_{\text{Shrinkage}}$	$\underbrace{W^3D + (P-1)W^2D}_{\text{Linear systems}} + PJ \cdot (\underbrace{WD \log(D)}_{\text{FFTs}} + \underbrace{WD}_{\text{Shrinkage}})$
Heide et al. [15]	$\underbrace{W^3D + (P-1)W^2D}_{\text{Linear systems}} + PJ \cdot (\underbrace{WD \log(D)}_{\text{FFTs}} + \underbrace{WD}_{\text{Shrinkage}})$	$\underbrace{W^3D + (P-1)W^2D}_{\text{Linear systems}} + PJ \cdot (\underbrace{WD \log(D)}_{\text{FFTs}} + \underbrace{WD}_{\text{Shrinkage}})$	$\underbrace{W^3D + (P-1)W^2D}_{\text{Linear systems}} + PJ \cdot (\underbrace{WD \log(D)}_{\text{FFTs}} + \underbrace{WD}_{\text{Shrinkage}})$
CCSC	$\underbrace{\frac{1}{U} (W^3D + (P-1)W^2D)}_{\text{Linear systems}}$	$\underbrace{\frac{1}{U} PJ \cdot (WD \log(\frac{D}{N}) + WD)}_{\text{FFTs}}$	$\underbrace{\frac{1}{U} PJ \cdot (WD \log(\frac{D}{N}) + WD)}_{\text{Shrinkage}}$

Dataset Size	CSC [15]	CCSC ($U = \# \text{PCs}$)		
		$U=5$	$U=10$	$U=50$
100	203.56 sec	35.35 sec	25.69 sec	25.20 sec
500	1530.71 sec	259.30 sec	82.69 sec	28.57 sec
1000	Out of Memory	387.68 sec	255.38 sec	35.63 sec

Figure 4: Complexity and Runtime Analysis. **Top:** Theoretical per-iteration cost of CCSC and other current CSC methods. **Bottom:** Runtime comparisons between the best competing CSC method [15] and CCSC. We demonstrate the runtime gain for a varying number of parallel working threads (U) and increasing dataset size.

parallel, affecting the runtime of the individual iterations. With full parallelization CCSC closely matches classical, non-distributed runtimes, while at the same time allowing CSC to scale. We first present the theoretical computational cost for a single iteration in Figure 4 (top), with P being the number of inner iterations (of the substeps in Alg. 2) and $U \leq N$ being the number of parallel workers. Assuming N blocks of equal size, splitting and distributing drastically reduces the cost of the linear system solves and of the fourier transforms. In terms of runtime, this smaller per-iteration cost allows more iterations in the proposed consensus optimization, while at the same time enabling scalability in terms of the memory requirements.

In Figure 4 (bottom) we provide empirical evidence of the high computational efficiency of the proposed approach by comparing the best competing CSC technique [15] with CCSC for increasing sizes of 2D dataset with varying number of parallel workers. For example, with a 2D dataset composed of 500 examples (each 100×100 pixels), we observe a speedup of $19\times$ for 10 workers, and a $54\times$ for 50 workers over existing CSC methods. Please note that, for datasets of larger size, current CSC techniques are intractable. All algorithms were executed on Intel Xeon 2.7 GHz Dual-core processor with 128GB RAM.

4. Learning

Large-scale Feature Learning on ImageNet: To test CCSC on large-scale image data, we use it to learn a dictionary for 5000 images from ImageNet [9] which is at least an order of magnitude more images than previously feasible with CSC methods. The dictionary itself consists of 100 filters of size 11×11 , and can be seen in Figure 1. For comparison we have included a similar dictionary trained on a very small fruit dataset. Although superficially similar, the large scale dictionary contains more general features which

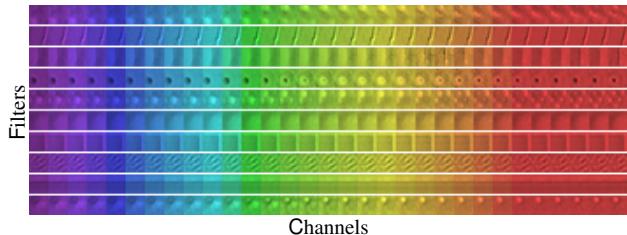


Figure 5: Multispectral (2D convolutional + 1D non-convolutional) dictionary. We show examples of 10 filters (out of a total of 100) learned across all 31 channels on the CAVE dataset. Note the similarity in the kernels across channels which depict the inherent redundancy along multi-channel image data.

lead to better reconstruction results (Sec. 5). Our dictionary also contains noise-like filters similar to those learned by discriminative feature learning models [8].

Multi-Spectral Feature Learning: Next, we test CCSC on multispectral data. Each image is now a 3-dimensional entity, with the wavelength as the extra dimension. However, this third dimension is typically much smaller (31 channels in our case) than the two spatial dimensions, and thus we chose to convolve only along the spatial dimensions while the third dimension is non-convolutional in the CCSC dictionary. We therefore force each pixel in the image to share the same coefficients for each element in the dictionary which promotes similarity among all channels without the need for any group sparsity constraints. We found that this method was greatly superior to solving each channel individually with 2D CSC, particularly in the presence of missing data where the proposed method is able to pull information across all channels. For details please refer to the supplementary material.

We trained the dictionaries on a select number of images from the Foster et al. [12] and CAVE [30] hyperspectral datasets, each learning 100, $11 \times 11 \times 31$ filters. An example of the CAVE filters can be seen in Figure 5 which show how the proposed framework learns a variety of features that slowly vary from channel to channel.

Video Feature Learning: Unlike multispectral data which contains a fixed number of channels, videos are composed of an arbitrary number of frames which lends itself to a fully convolutional 3D filter. Therefore, we learned a set of 49 3D filters of size $11 \times 11 \times 11$ from a varied set of 64 HD video clips. A sample set of these filters can be seen in Figure 6, which demonstrates the variety of CCSC filters as well as their smooth spatial and color transitions across time frames. For reconstruction results please refer to the supplemental material.

Light Field Feature Learning: Although typically captured as a single image, light fields can be represented as a 4D tensor with two spatial dimensions and two angular dimensions. Because the two angular dimensions are

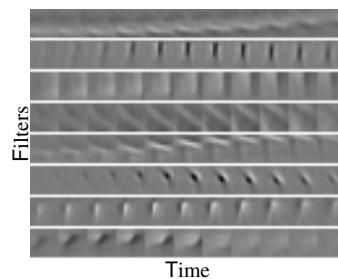


Figure 6: Learned Video Features (3D-Convolutional). Each row shows a single 3D convolutional video kernel whose features slowly change over time from left to right.

small (typically only 5 to 8 angles), we chose to train dictionary filters which were convolutional spatially, but non-convolutional in the angular dimensions. The final dictionary was trained on a set of 64 light fields truncated to 5 angular views in both x and y, and contained 49 filters of size $11 \times 11 \times 5 \times 5$. A sample set of these filters can be found in Figure 7 which clearly demonstrates the angular structure learned by CCSC. Each 5×5 group of filters slowly varies across the angular dimensions while exhibiting general features for reconstruction throughout.

5. Reconstruction

M-Operator: Similar to Heide et al. [15], we employ a binary mask M as a general linear operator which can be used for a variety of purposes, such as boundary handling, and masking incomplete data. Note that, typically M is a diagonal or block diagonal matrix, such that it decouples linear systems of the form $(M^T M + I)x = v$ into many small independent systems that can be efficiently solved.

Inpainting and Deconvolution: To compare the CCSC large-scale dictionary with conventional CSC, and demonstrate applicability to different noise and image formation models, we evaluated their performance in both inpainting and Poisson noise deconvolution with the Poisson proximal

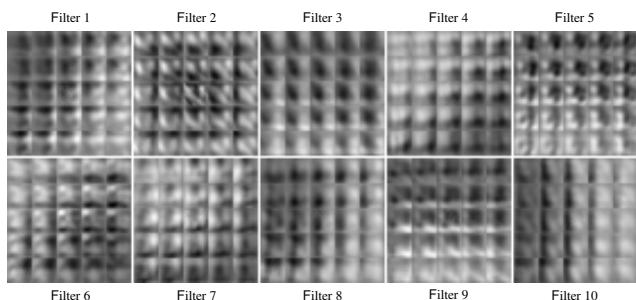


Figure 7: Example of 10 Learned Light Field Features (2D Convolutional + 2D Non-Convolutional). Each group of 5×5 filters shows all 25 angular features learned.

648
649
650
651
652
653
654
655
656
657

2D Inpainting

Image	CCSC	CSC [15]	NLR [10]
Wind Mill	35.13	33.49	27.30
Sea Rock	28.45	27.29	23.38
Parthenon	31.36	29.79	24.99
Rolls Royce	29.15	27.34	22.05
Fence	30.83	29.59	22.41
Car	34.06	32.57	22.86
Kid	29.41	28.05	23.08
Tower	29.97	28.30	24.86
Fish	31.68	30.26	20.92
Food	36.77	35.09	23.78

2D Poisson Deconvolution

Image	CCSC	CSC [15]	Krishnan [19]
Agama	28.05	27.87	24.26
Gyful	29.36	29.31	24.52
Kathmandu	23.14	22.86	20.19
Laser	31.59	31.57	28.10
Libelle	27.56	27.17	22.86
Melinaea	28.36	28.05	24.16
Mototaxis	23.41	23.34	20.87
Painted	22.47	21.89	18.84
Platycercus	25.93	24.94	21.12
Porsche	27.11	26.21	19.60

3D Multispectral Demosaicing

Image	CCSC	IID [24]	SD [4]	WB [4]
Balloons	28.62	27.38	25.83	25.91
Beads	23.87	23.33	14.18	14.57
CD	30.54	23.38	18.23	18.39
Chart	24.64	21.56	13.92	13.84
Clay	29.74	14.25	12.53	12.57
Cloth	23.50	20.96	13.91	14.02
Statue	33.38	20.83	16.97	17.12
Face	28.12	17.50	12.91	13.00
Beer	23.72	18.33	9.21	9.23
Food	28.91	25.50	17.38	17.61

Figure 8: Quantitative analysis of 2D Image Reconstruction and Multispectral Demosaicing. **Left:** Inpainting results for 50% randomly subsampled observations of images randomly selected from ImageNet [9]. The filters learned using CCSC (shown in Fig. 1) lead to significantly prediction results compared to the ones from [15], as well as recent patch-based methods such as the non-local low-rank method from [10]. **Center:** 2D Poisson Deconvolution. Comparisons of CCSC against the state of the art deconvolution method [19] and the classical CSC method. **Right:** Multispectral Demosaicing results for the CAVE dataset comparing CCSC against the state of the art Iterative Intensity Difference (IID) [24], and the previous standard Spectral Difference (SD) [4] and Weighted Bilinear (WB) [4] interpolation methods. All values reported as PSNR in dB. Please see supplement for comparisons of CCSC with other state of the art techniques.

operator described in the supplement. Quantitative results can be found in Figure 8 (left and center), and sample reconstructions can be found in Figure 9 & 10. In all cases the CCSC large-scale features outperformed both classical CSC as well as state of the art alternatives. Please see supplement for additional comparisons of our algorithm with other state of the art techniques.

Multi-Spectral Demosaicing: We compare the proposed method to the state of the art multispectral demosaicing technique [24]. To emulate the demosaicing process we process the raw data to conform to a multispectral filter array (MSFA) pattern with 16 evenly spaced channels corresponding to data from the 400 to 700 nm range. We then reconstruct the data as a sub-sampling problem where the missing data from each channel is masked by the M operator. We compared the CCSC results with the code provided by [24] on the original CAVE dataset [30] and calculated the PSNR of the entire reconstructed image. The results in Figure 8 (right) show that CCSC outperforms state of the art techniques, an example of which can be seen in Figure 11.

Light Field View Synthesis: Here we compare CCSC using the learned light field dictionary with state of the art light field view synthesis algorithms. The results can be found in Figure 12 along with sample output. Using the M operator to mask the unknown views we wish to synthesize, we can employ our general reconstruction algorithm to generate the missing data. Using the dictionary described in previous sections with 5×5 angular views and testing data provided by [16], we synthesized the second and fourth angular views in both x and y after removing them from the data. Although this is not the experimental setup used in [16], which may account for some degradation in their performance, it demonstrates the versatility of the proposed approach. One dictionary trained with CCSC can be used to

synthesize any number or orientation of light field views.

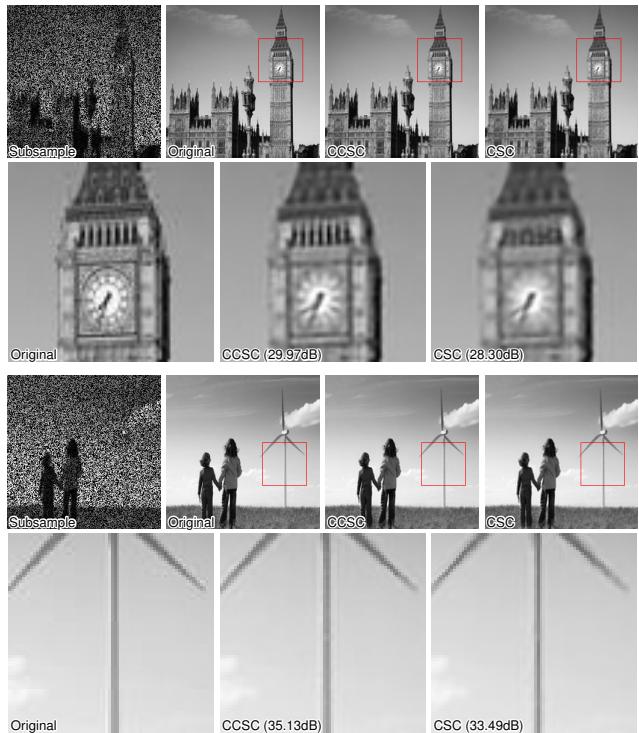


Figure 9: Inpainting results using 2D filters for both “Clock” and “Wind Mill”. Top rows show from left to right: (a) Subsampled image, (b) Ground Truth, (c) CCSC, (d) Conventional CSC. Bottom shows insets from (b-d) respectively. It is evident that conventional filters fail for difficult contrast edges (vertical clock features) and overfit for other structures (white region of wind turbine stand).



Figure 10: Deconvolution results using 2D filters for both Bird and Car. Top rows show from left to right: (a) Blurred image, (b) Ground Truth, (c) CCSC, (d) Conventional CSC. Bottom shows insets from (b-d) respectively. In both darker regions of bird feathers and car text conventional CSC hallucinates features which are not present resulting in poor deconvolution results.

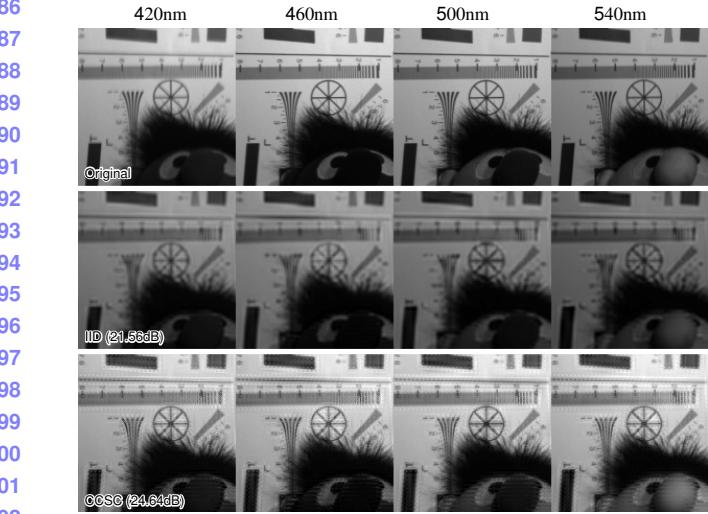


Figure 11: Multispectral demosaicing results from four wavelengths of the chart dataset. Note that while the proposed algorithm does contain demosaicing artifacts it is better able to reconstruct the high frequency details found in the background chart while preserving spectral differences. Find a comparison to additional methods, WB (13.84 dB) and SD (13.92 dB), in the supplement.

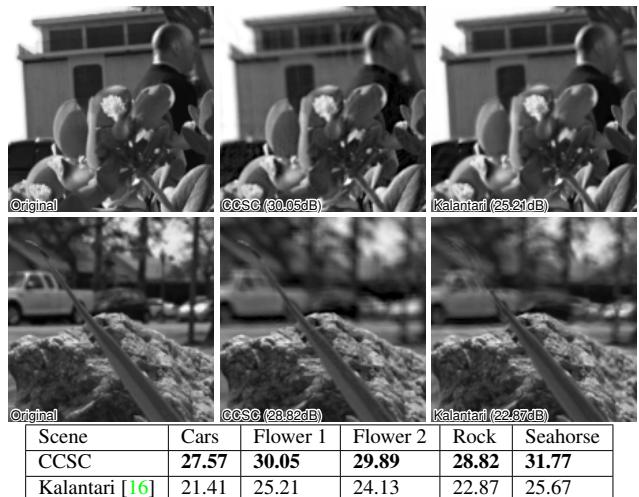


Figure 12: **Top:** Example of two synthesized views from the Flower 1 and Rock datasets. From left to right, (a) Ground Truth, (b) CCSC, (c) Kalantari [16]. The proposed algorithm produces less noticeable ghosting artifacts due to far away objects and better reconstructs fine detail in nearby objects such as the leaf edges and stalk tip. **Bottom:** Quantitative reconstruction results in PSNR (dB).

6. Discussion

Conclusion We have shown that CSC has the potential to be applied in many high and low level computer vision applications. Our distributed CCSC algorithm is both memory efficient and capable of high quality representations of N-Dimensional image data. Furthermore, by reducing and distributing the memory requirements compared to previous CSC methods, our algorithm is capable of handling much larger datasets thereby generating more generalized feature spaces. With our proposed method, we hope to provide a step towards practical and efficient approaches to solving high-dimensional sparse coding problems.

Future Work Although we have shown that CCSC is capable of tackling many computer vision problems, there are many further possible applications. Because our algorithms produce high-dimensional per-pixel coefficients, they could be incorporated into classification, segmentation, or spectral unmixing techniques.

Unlike previous CSC implementations, our distributed framework is amenable to GPU implementation which often have extreme memory constraints. Such an implementation would dramatically increase performance and, for example, bring our multispectral demosaicing algorithm run time in line with other methods.

References

- [1] M. Aharon, M. Elad, and A. Bruckstein. K-svd: An algorithm for designing overcomplete dictionaries for

- 864 sparse representation. *IEEE Trans. Signal Processing*,
865 54(11):4311–4322, 2006. 1, 5
- 866 [2] N. Akhtar, F. Shafait, and A. Mian. Bayesian sparse
867 representation for hyperspectral image super resolution.
868 In *Proc. IEEE CVPR*, pages 3631–3640, 2015.
869 1
- 870 [3] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eck-
871 stein. Distributed optimization and statistical learning
872 via the alternating direction method of multipliers.
873 *Foundations and Trends in Machine Learning*, 3(1):1–
874 122, 2011. 3, 4
- 875 [4] J. Brauers and T. Aach. A color filter array based
876 multispectral camera. In *12. Workshop Farbbildver-
877 arbeitung*. Ilmenau, 2006. 7
- 878 [5] H. Bristow, A. Eriksson, and S. Lucey. Fast convolu-
879 tional sparse coding. In *Proc. CVPR*, pages 391–398,
880 2013. 1, 3, 4, 5
- 881 [6] H. Bristow and S. Lucey. Optimization methods for
882 convolutional sparse coding. *arXiv:1406.2407*, 2014.
883 1, 3, 4, 5
- 884 [7] A. M. Bruckstein, D. L. Donoho, and M. Elad. From
885 sparse solutions of systems of equations to sparse
886 modeling of signals and images. *SIAM review*,
887 51(1):34–81, 2009. 1
- 888 [8] Y. Chen, W. Yu, and T. Pock. On learning opti-
889 mized reaction diffusion processes for effective image
890 restoration. In *Proc. IEEE CVPR*, pages 5261–5269,
891 2015. 6
- 892 [9] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and
893 L. Fei-Fei. Imagenet: A large-scale hierarchical im-
894 age database. In *Proc. IEEE CVPR*, pages 248–255,
895 2009. 1, 5, 7
- 896 [10] W. Dong, G. Shi, X. Li, Y. Ma, and F. Huang.
897 Compressive sensing via nonlocal low-rank regulariza-
898 tion. *IEEE Trans. Image Processing*, 23(8):3618–
899 3632, 2014. 7
- 900 [11] M. Elad and M. Aharon. Image denoising via sparse
901 and redundant representations over learned dictionar-
902 ries. *IEEE Trans. Image Processing*, 15(12):3736–
903 3745, 2006. 1
- 904 [12] D. H. Foster, K. Amano, S. M. Nascimento, and M. J.
905 Foster. Frequency of metamerism in natural scenes.
906 *JOSA A*, 23(10):2359–2372, 2006. 6
- 907 [13] A. Graves, A.-r. Mohamed, and G. Hinton. Speech
908 recognition with deep recurrent neural networks. In
909 *Proc. IEEE ASSP*, pages 6645–6649. IEEE, 2013. 1
- 910 [14] R. B. Grosse, R. Raina, H. Kwong, and A. Y. Ng.
911 Shift-invariance sparse coding for audio classifica-
912 tion. In *Proc. UAI*, pages 149–158, 2007. 1
- 913 [15] F. Heide, W. Heidrich, and G. Wetzstein. Fast and
914 flexible convolutional sparse coding. In *Proc. IEEE
915 CVPR*, pages 5135–5143, 2015. 1, 2, 3, 4, 5, 6, 7
- 916 [16] N. K. Kalantari, T.-C. Wang, and R. Ramamoorthi.
917 Learning-based view synthesis for light field cameras.
918 *arXiv preprint arXiv:1609.02974*, 2016. 7, 8
- 919 [17] K. Kavukcuoglu, P. Sermanet, Y. Boureau, K. Gre-
920 gor, M. Mathieu, and Y. LeCun. Learning convo-
921 lutional feature hierarchies for visual recognition. In
922 *Proc. NIPS*, 2010. 1
- 923 [18] B. Kong and C. C. Fowlkes. Fast Convolutional Sparse
924 Coding (FCSC). Technical report, UCI, May 2014. 3
- 925 [19] D. Krishnan and R. Fergus. Fast image deconvolu-
926 tion using hyper-laplacian priors. In *Proc. NIPS*, pages
927 1033–1041, 2009. 7
- 928 [20] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Im-
929 agenet classification with deep convolutional neural
930 networks. In *Proc. NIPS*, pages 1097–1105, 2012. 1
- 931 [21] X. Lin, Y. Liu, J. Wu, and Q. Dai. Spatial-spectral
932 encoded compressive hyperspectral imaging. *ACM
933 Trans. Graphics*, 33(6):233, 2014. 1
- 934 [22] J. Mairal, F. Bach, J. Ponce, and G. Sapiro. Online
935 dictionary learning for sparse coding. In *Proc. ICML*,
936 pages 689–696. ACM, 2009. 1
- 937 [23] K. Marwah, G. Wetzstein, Y. Bando, and R. Raskar.
938 Compressive light field photography using overcom-
939 plete dictionaries and optimized projections. *ACM
940 Trans. Graph. (SIGGRAPH)*, 32(4):46:1–46:12, 2013.
941 1
- 942 [24] S. Mihoubi, O. Lossen, B. Mathon, and L. Macaire.
943 Multispectral demosaicing using intensity-based spec-
944 tral correlation. In *Proc. IEEE IPTA*, pages 461–466,
945 2015. 7
- 946 [25] B. A. Olshausen and D. J. Field. Sparse coding with
947 an overcomplete basis set: A strategy employed by
948 v1? *Vision Research*, 37(23):3311 – 3325, 1997. 1
- 949 [26] N. Parikh and S. Boyd. Proximal algorithms. *Founda-
950 tions and Trends in Optimization*, 1(3):123–231, 2013.
951 4
- 952 [27] P. Sermanet, D. Eigen, X. Zhang, M. Mathieu, R. Fer-
953 gus, and Y. LeCun. Overfeat: Integrated recogni-
954 tion, localization and detection using convolutional
955 networks. *arXiv preprint arXiv:1312.6229*, 2013. 1
- 956 [28] K. Simonyan and A. Zisserman. Very deep convo-
957 lutional networks for large-scale image recognition.
958 *arXiv preprint arXiv:1409.1556*, 2014. 1
- 959 [29] J. Yang, J. Wright, T. S. Huang, and Y. Ma. Im-
960 age super-resolution via sparse representation. *IEEE
961 Trans. Image Processing*, 19(11):2861–2873, 2010. 1
- 962 963 964 965 966 967 968 969 970 971

- 972 [30] F. Yasuma, T. Mitsunaga, D. Iso, and S. K. Nayar. 1026
973 Generalized assorted pixel camera: postcapture 1027
974 control of resolution, dynamic range, and spectrum. *IEEE 1028
975 Trans. Image Processing*, 19(9):2241–2253, 2010. 6, 1029
976 7 1030
977 [31] M. D. Zeiler, D. Krishnan, G. W. Taylor, and R. Fer- 1031
978 gus. Deconvolutional networks. In *Proc. CVPR*, pages 1032
979 2528–2535, 2010. 1, 5 1033
980 [32] M. D. Zeiler, G. W. Taylor, and R. Fergus. Adaptive 1034
981 deconvolutional networks for mid and high level fea- 1035
982 ture learning. In *Proc. ICCV*, pages 2018–2025, 2011. 1036
983 1 1037
984 [33] D. Zoran and Y. Weiss. From learning models of natu- 1038
985 ral image patches to whole image restoration. In *Proc. 1039
986 ICCV*, pages 479–486, 2011. 1 1040
987 1041
988 1042
989 1043
990 1044
991 1045
992 1046
993 1047
994 1048
995 1049
996 1050
997 1051
998 1052
999 1053
1000 1054
1001 1055
1002 1056
1003 1057
1004 1058
1005 1059
1006 1060
1007 1061
1008 1062
1009 1063
1010 1064
1011 1065
1012 1066
1013 1067
1014 1068
1015 1069
1016 1070
1017 1071
1018 1072
1019 1073
1020 1074
1021 1075
1022 1076
1023 1077
1024 1078
1025 1079