# Deep learning based dictionary learning and tomographic image reconstruction

**Jevgenija Rudzusika**
Department of Mathematics
KTH Royal Institute of Technology
jevaks@kth.se

**Thomas Koehler**
Philips Research
thomas.koehler@philips.com

**Ozan Öktem**
Department of Mathematics
KTH Royal Institute of Technology
ozan@kth.se

**Abstract**

This work presents an approach for image reconstruction in clinical low-dose tomography that combines principles from sparse signal processing with ideas from deep learning. First, we describe sparse signal representation in terms of dictionaries from a statistical perspective and interpret dictionary learning as a process of aligning distribution that arises from a generative model with empirical distribution of true signals. As a result we can see that sparse coding with learned dictionaries resembles a specific variational autoencoder, where the decoder is a linear function and the encoder is a sparse coding algorithm. Next, we show that dictionary learning can also benefit from computational advancements introduced in the context of deep learning, such as parallelism and as stochastic optimization. Finally, we show that regularization by dictionaries achieves competitive performance in computed tomography (CT) reconstruction comparing to state-of-the-art model based and data driven approaches.

## 1 Introduction

This work presents an approach for image reconstruction in clinical low-dose tomography that combines principles from sparse signal processing with ideas from deep learning. Before describing the method, we begin with providing some background in this section that serves as a motivation for our work.

### 1.1 X-ray tomographic imaging in medicine

In CT an x-ray tube and a detector rotate around a patient, acquiring x-ray transmission measurements from multiple directions. The aim is to computationally recover an image showing the interior anatomy of the patient. This is an inherently unstable procedure, so a key issue lies in adequately addressing this instability.

CT is nowadays one of the most frequently used imaging modalities and approximately 100 million CT scans are being performed annually in the United States alone. The repeated radiation from multiple directions means exposing

the patient to ionising x-rays. Rising concerns about radiation dose in clinical CT imaging lead to increasing interest in low-dose protocols [18, 68, 63]. One approach to lower the dose is to collect data along fewer directions (sparse-view CT). Alternatively, one can keep full sampling of data and instead reduce the intensity of emitted x-ray photons, e.g., by lowering the tube current and/or voltage in the scanner, which in turn results in data that is significantly more noisy. This second approach is not only more practical, it is also more efficient, i.e., one can obtain better image quality from data with the same relative dose reduction [36, 50]. This is true at least for the standard baseline methods. Therefore, all of our low-dose CT data refers to data from full angular sampling, but with high noise.

The significance of how the aforementioned instability is addressed during reconstruction increases with the noise level in data. In particular, it is evident that the computationally efficient reconstruction methods originally developed for image reconstruction from less noisy normal-dose CT data are inadequate for low-dose CT. The resulting images are degraded by noise and artefacts, which in turn renders them sub-optimal for diagnostic interpretation. Addressing this has been an active area for research as briefly outlined in section 2.1.

## 1.2   Inverse problems and regularization

CT imaging is an example of an inverse problem. The latter refers to problems where the goal is to recover a hidden signal from measured data that represent noisy indirect observations of the signal. Such problems arise in many areas of science and engineering.

Many inverse problems are large-scale in the sense that data and/or signal reside in high-dimensional spaces even after clever discretization. As an example, signals and data in 2D/3D tomographic imaging are represented by high dimensional arrays. Solution methods must therefore have computationally feasible implementations, which is especially important in time critical applications.

A key challenge in solving inverse problems is to address *instability (ill-posedness)*. This refers to a situation where merely maximising the fit against measured data is an unstable procedure, i.e., small errors in data result in large perturbations to the resulting signal. Solving an ill-posed inverse problem therefore requires specific attention to handle the intrinsic instability. *Regularization* refers to mathematical theory and algorithms that introduce stability by balancing the need to fit data against having a reconstruction that is consistent with known information about the unknown signal one seeks to recover (prior model). The data misfit can be quantified by the *data (log) likelihood*, which consists of a *forward operator* (models how a signal gives rise to data in absence of noise) and a *noise model* (encodes statistical properties of the observation errors). Both these constitute a simulator and they are typically derived from first principles by carefully considering the physics that governs the formation of data. Prior models are, in contrast to the data likelihood, not derived from the physics of data acquisition. Next, the impact from a specific choice of prior becomes increasingly notable as the noise level in data increases. Hence, a major topic during the last two decades in inverse problems research has been to devise appropriate prior models. Early approaches had smoothing priors for recovering low-frequency components of the signal. These were followed by more intricate regularity priors, like more domain adapted sparsity promoting priors that are

either handcrafted or learned from example data.

Overall, it is still challenging to suggest a reconstruction method for low-dose CT that is *computational feasible*, yet is based on a *prior model* that provides a clinically notable improvement in image quality.

## 1.3  Priors given by generative models

A common approach for assembling a domain adapted prior model is to start from a *generative model* that is defined by a *synthesis operator* $\mathcal{S}\colon Z \to X$, which generates a signal in $X$ from a representation in *latent space* $Z$. This representation is often not unique, therefore a criterion for choosing most suitable representation is introduced. This criterion can usually be expressed as minimising a (regularization) functional $\mathcal{R}\colon Z \to \mathbb{R}$. The generative model and the aforementioned regularization functional define a prior model, which can be used to address the inherit instability in solving an ill-posed inverse problem as outlined in section 5.

**Dictionaries and sparsity**   A sparsity prior assumes that signals from some vector space $X$ can be described as a *linear combinations* of a *few* signal components (atoms) from a pre-specified *dictionary*.

Stated formally, a *dictionary* is some countable subset $\mathcal{D} = \{d_i\}_i \subset X$ whose elements $d_i$ are called *(dictionary) atoms*. The synthesis operator defining the generative model assembles signals in $X$ by taking linear combinations of atoms:

$$\mathcal{S}_{\mathcal{D}}(z) = \sum_i z_i d_i \quad \text{for } z_i \in \mathbb{R} \text{ and } d_i \in \mathcal{D} \subset X. \tag{1}$$

The latent variables $z_i \in \mathbb{R}$ are called *dictionary coefficients*.

A sparsity prior implies that only few non-zero coefficients are sufficient to represent a signal. Unfortunately, the generative model in eq. (1) in combination with a sparsity prior would require too many dictionary atoms to represent large and complex signals like images. Therefore, it is typically applied to smaller image patches and representation of the full image is obtained by averaging representations of overlapping patches. An alternative is to consider convolutional dictionaries:

$$\mathcal{S}_{\mathcal{D}}(z) = \sum_i z_i * d_i \quad \text{for } z_i \in X \text{ and } d_i \in \mathcal{D} \subset X. \tag{2}$$

The dictionary atoms $d_i$ represent small local features, i.e. have a relatively small support compared to signals in $X$. They synthesise a signal by convolving against (local) signal maps, so dictionary coefficients $z_i$ are elements in $X$ (coefficient maps) instead of simply real numbers.

The problem of computing sparse dictionary coefficients for a given signal $x$ is called sparse coding. It is natural to define this problem as finding coefficients $\widehat{z}$ that minimize the objective:

$$\min_{z \in Z} \lVert \mathcal{S}_{\mathcal{D}}(z) - x \rVert_X^2 + \lambda \lVert z \rVert_0. \tag{3}$$

The $l_0$ pseudo-norm in $\lVert z \rVert_0$ simply counts the number of non-zero terms in the sequence $z$. Since each non-zero term corresponds to an atom in $\mathcal{D}$, minimizing

the objective means we look for a sparse representation (the one with few number of atoms), which sufficiently well approximates the original signal. In particular, $\lambda \geq 0$ is a weighting factor between sparsity and how well the coding matches the signal. A small value means the coding has to be close to the signal at the expense of sparsity, whereas increasing $\lambda$ means prioritising sparsity over the need to match the signal.

Unfortunately, the presence of an $l_0$-pseudo norm in eq. (3) makes the optimization problem hard to solve. Indeed, it is known that this problem is NP-hard [5, 69]. This difficulty is addressed by using approximation algorithms. There are generally two types of approximation techniques for this purpose. Greedy algorithms, like orthogonal matching pursuit [71] and iterative hard thresholding [14], address the problem eq. (3) directly, while methods based on convex relaxation replace the $l_0$-pseudo norm with an $l_1$-norm [21]:

$$\min_{z \in Z} \left\| \mathcal{S}_{\mathcal{D}}(z) - x \right\|_X^2 + \lambda \|z\|_1. \tag{4}$$

In general, a solution to eq. (3) yields representations that can be more sparse than the representation obtained from solving eq. (4). However, if $x$ admits a sufficiently sparse representation in $\mathcal{D}$ and $\mathcal{D}$ satisfies the restricted isometry property, then a solution to eq. (4) is with high probability also a solution to eq. (3) [25, 19, 20].

**Dictionary learning**  A good sparsifying dictionary will generate sparse codes of signals that preserve important features, whereas noise and artefacts are preferably suppressed (or even lost). The choice of dictionary is therefore an essential component, and in general, this can be done using one of two ways: (a) building a dictionary based on a mathematical model of signals in $X$, or (b) learning a dictionary to perform best on a training set [61]. Dictionary learning focuses on the latter, i.e., on the task of learning an appropriate sparse representation from example data. This is an inverse problem whose solution is a (trained) dictionary, which in turn defines a prior model on $X$.

Let $x_1, \ldots, x_n \in X$ be the example signals and $\mathscr{D}$ is some fixed family of dictionaries (countable subsets of $X$). A common approach is to perform dictionary learning jointly with sparse coding, e.g., by minimizing the following objective with respect to both, dictionaries $\mathcal{D} \in \mathscr{D}$ and coefficients $z_i \in Z$:

$$\min_{\mathcal{D}, z_i} \sum_{i=1}^{n} \left\| \mathcal{S}_{\mathcal{D}}(z_i) - x_i \right\|_X^2 + \lambda \|z_i\|_1, \tag{5}$$
$$\text{s.t.} \ \|d_j\| = 1 \quad \text{for all } d_j \in \mathcal{D}.$$

The constraint on the norm of dictionary atoms $\|d_j\| = 1$ is included to circumvent the fact that coefficients $z_i$ can be reduced simply by up-scaling the corresponding dictionaries.

In section 3 we present a statistical formulation of dictionary representation that provides a mathematical interpretation of the dictionary learning scheme in eq. (5).

**Dictionary based CT reconstruction**  Our goal here is to solve the inverse problem arising in CT imaging (*image reconstruction*) assuming the true (unknown)

image is generated by the aforementioned generative model. This assumption may serve as a regularization, which is in particular the case when the generative model is not capable of generating undesirable noise and artefacts.

These two inverse problems, i.e., defining the generative model and image reconstruction, can be solved sequentially or jointly. The former means that dictionary learning is performed without considering the CT image reconstruction step that follows next.

In this work, we try to benefit from computational advancements of deep learning while staying in a highly interpretative framework of representation through dictionaries. In this respect our work is strongly related to [70]. However, our main focus is on dictionary learning as a regularization method in inverse problems, specifically in CT.

## 1.4 Outline of paper

In the next section (section 2) we present a survey of related works on dictionary learning and reconstruction in computerized tomography. Section 3 presents dictionary learning from a statistical perspective as a way to approximate an empirical distribution of a natural signal. Secondly, section 3.3 describes a practical implementation of the learning procedure. This is based on the convolutional generative model in eq. (2) and the sparse coding in eq. (4) that defines a prior model parametrised by a learned dictionary. Next, in section 4 we describe how the learned dictionary model can be used for regularization in inverse problems. In particular, section 5 describes our reconstruction experiments in computerized tomography. Section 6 concludes the paper.

## 2 Survey of related work

### 2.1 Image reconstruction in clinical low-dose CT

Image reconstruction for clinical CT has traditionally relied on the filtered back-projection (FBP) method and variants thereof. These analytical methods were first in introduced in late 1960s in astronomy [15] and later adopted by the medical imaging community in the early 1970s [64]. They use principles from Fourier analysis and sampling theory to recover the part of the image that is band-limited whereas high frequency components, like noise, are filtered out. The mathematical foundation of FBP was developed in the late 1970s [52] and it has since then continuously evolved to account for increasingly complex acquisition geometries, like those that arise in 3D clinical CT [30, 29, 39, 55, 42, 78, 38, 75]. The compromise that FBP type of methods make in balancing image quality against reconstruction speed is still difficult to outperform in the context of clinical normal dose CT imaging [58].

FBP type of methods do not handle the noise statistics of measured data optimally. They regularise by recovering the band-limited part of the image, an approach that is not sufficient to suppress noise and artefacts that degrade image quality in low-dose CT. Hence, FBP type of methods render images in low-dose CT that are sub-optimal for diagnostic interpretation. The need to address this issue has catalysed the development of new iterative reconstruction algorithms in both academia and industry. Only a fraction of the methods developed in

academia for low-dose CT have made it into the clinic practice, and this almost always as part of collaboration with vendors of CT scanners as surveyed next.

The first commercially available iterative reconstruction algorithms for replacing FBP appeared in 2009 with the launch of IRIS (Siemens Healthineers) and ASiR (GE Healthcare). Since then, within a few years, all major CT vendors introduced iterative reconstruction algorithms for clinical routine. Examples are AIDR (Canon Medical Systems, 2010) and AIDR3D (Canon Medical Systems, 2012), iDose$^4$ (Philips Healthcare, 2012), SAFIRE (Siemens Healthineers, 2011) and ADMIRE (Siemens Healthineers, 2014). These methods (often referred to as hybrid techniques or statistical iterative reconstruction) define an iterative scheme that combines denoising data and/or image with FBP to map data into image space. Variational models (also called model-based iterative reconstruction) were introduced somewhat later. These come with rigorous mathematical foundations and in addition to noise and photon statistics, one can also model object, scanner geometry and detector response. These are better at reducing noise and artefacts than hybrid techniques, but they may also alter image texture more. Examples are VEO (GE Healthcare, 2011), ASiR-V (GE Healthcare, 2015), IMR (Philips Healthcare, 2015), and FIRST (Canon Medical Systems, 2016). See [31, 6, 51, 73] for a review of these methods in low-dose setting.

The final line of development is the recent commercial usage of deep learning based approaches for denoising/image reconstruction, like TrueFidelity (GE Healthcare, 2018) and AiCE (Canon Medical Systems, 2018).

## 2.2 Dictionary learning

Traditional dictionary learning typically uses a variational model with a sparsity promoting regulariser, as in eq. (5), to recover the dictionary and corresponding dictionary coefficients for signals in a given training data set.

For a long time K-SVD algorithm [4] was among the best approaches for dictionary learning. This method updates dictionary atoms one by one, while using all the available training data. Such an algorithm is therefore not suitable for large training data-sets that are available nowadays. Convex relaxation is another, theoretically appealing approach for learning convolutional dictionaries. It was proposed in [22] and the advantage is that a unique solution to the relaxed problem can be determined. Unfortunately, it comes at a cost of greatly increasing the number of optimization variables, so such an approach is suitable only for learning a few relatively smalls atoms.

**Statistical formulation**   An alternative viewpoint is to phrase dictionary learning as density estimation. Dictionaries are here used to construct an estimate from training data of the unobservable probability density of signals.

Such an approach taken in eq. (7), which uses variational inference to minimize the Kullback-Leibler (KL) divergence between the empirical distribution of signals in the training data and the distribution modeled by the dictionaries. As shown in eq. (15), this naturally leads to a formulation where the joint probability density of signals and dictionary coefficients, $\rho_{x,z}(x, z)$, is marginalized over the dictionary coefficients. This is computationally intractable for imaging applications. One option is to approximate the $Z$-integral in eq. (15) by the value of the integrand at the mode [57, 70], which in turn leads to the common

formulation for dictionary learning given in eq. (28). However, the above cited works do not characterize the resulting approximation error.

Another option is to approximate the integrand in eq. (15) by un-normalized Gaussian density around the mode [44, 45]. This is equivalent to approximating the posterior density $\rho_{z|x}(z \mid x)$ by a Gaussian distribution and it leads to an optimization problem different than the one in eq. (28). Yet another dictionary learning algorithm is derived in [32]. Here, the authors start out from expressing the Laplace prior $\rho_z(z)$ in eq. (15) as a supremum of Gaussian densities. Then, they optimize a lower bound of the log-likelihood. Unfortunately, the derivation of the algorithms in these cited papers involve inversion of prohibitively large matrices that require further approximations for their proper handling.

Our aim with section 3 is not primarily to derive a new update rule for dictionary learning as in [44, 45, 32]. It is rather to provide a statistical interpretation of eq. (28) that, unlike [57, 70], also includes characterizing the approximation error.

**Connection to deep learning**   There are many attempts at connecting sparse signal processing to deep learning. A popular line of investigation focuses on the connection between sparse representation through convolutional dictionaries and convolutional neural networks.

More precisely, the authors of [59] show that a trained convolutional neural network corresponds to an approximate algorithm for the sparse coding in a sparse multi-layer convolutional dictionary model. This model exploits a structure given by an ordered sequence hierarchically arranged dictionaries. Starting from the signal, atoms in one dictionary are sparsely represented by atoms in the next dictionary. Thus, each subsequent dictionary layer represents features with higher level of abstraction and [66, 67, 48] proposes training strategies for such models. The obtained sparse representations with respect to the last dictionary layer are used as input for image classification.

In [48], it is empirically demonstrated that classification results based on such a trained sparse convolutional dictionary model are less susceptible to adversarial image perturbations than corresponding outcomes from a trained convolutional neural network.

## 2.3   Dictionary based CT reconstruction

Several authors attempted regularization by dictionary learning in CT. Early work used dictionary representation of image patches as a regularizing component in statistical iterative reconstruction (SIR) [74]. The authors evaluated a dictionary learned from training images (GDSIR) and an adaptive dictionary learned simultaneously with reconstruction (ADSIR) and found that both are similarly effective. GDSIR and ADSIR set the baseline for many works to follow. In this work we use a similar problem formulation. However, the way we address practical aspects of dictionary learning and sparse coding is different.

Later, in [24] one observed that artifacts in low-dose CT images might give rise to sufficiently large coefficients in dictionary representation and those high coefficients won't be suppressed by thresholding operation, unless a very high threshold is set and an image becomes over-smoothed. They use carefully selected samples to train dictionaries that represent artifacts as well as tissue features and then set coefficients corresponding to artifacts to 0. They evaluated their method

on real data and report benefits of their approach in qualitative assessment performed by radiologists. Nevertheless, the approach falls into category of post-processing reconstructions obtained by FBP [35, 53], while we perform reconstruction and de-noising jointly by solving a variational problem.

Other attempts to improve the performance of dictionary learning include: using $l_1$ for misfit between image and its dictionary representation [76], smoothing intermediate image updates to remove artifacts [43], using dictionary learning in combination with total variation (TV) [77] and clustering patches and learning dictionaries for each class separately [37].

Following the success of convolutional filters in deep learning, regularization by learned convolutional dictionaries has been applied to CT. In [10], one learns 32 convolutional dictionaries of size $10 \times 10$, while using total variation type penalty on coefficient maps.

Furthermore, most of the previously published work use relatively small dictionaries. A common setup is to learn 256 dictionaries of size $8 \times 8$ with 5 to 10 non-zero coefficients as it was done in [74]. One exception is the work in [33], where authors experimented with different patch sizes and found that the optimal patch size is around $16 \times 16$, with larger patches leading to slightly worse performance. In this work we also aim to learn larger dictionaries that can capture more information and potentially provide stronger regularization.

## 2.4   Deep Learning methods for CT image reconstruction

Most recently, the idea of non-linear synthesis operators was introduced. In [56], two networks, one for encoding and one for decoding, were trained, forcing encoders output to be sparse. Then, the decoder was used as a non-linear synthesis operator to solve regularized optimization problem. Similarly, in [8] convolutional neural network was used to generate images. Network parameters were trained in unsupervised manner, during the reconstruction itself. This approach called Deep Image Prior [72] is based on the idea that convolutional neural networks learn to represent low frequency components first and therefore, early stopping can be used to avoid noise. In contrast, we stick to the more traditional approach having linear synthesis operator.

Despite the lack of inseparability and theoretical guarantees, supervised deep learning methods have been shown to produce state of the art results in CT reconstruction. Although all of these methods share the philosophy of learning from data, there is a lot of diversity in the proposed model architectures and training procedures. These methods include Adversarial Regularizer [47], U-Net post-processing [34], and networks for reconstruction that are trained in a supervised manner with architectures derived from unrolling iterations of a suitable optimization scheme: LEARN [23], Learned Primal-Dual [2], and Total Deep Variation [41]. We refer to [60] for a survey of the recent methods.

## 3   Statistical interpretation of dictionary learning

To formulate dictionary learning in a statistical setting, we start by considering a *generative model* for signals in $X$. Let $x \sim P_X$ be a $X$-valued random variable generating natural signal, e.g., 2D/3D images. The generative model parametrised by dictionary $\mathcal{D} \subset X$ is defined as a signal generated by sampling

from the $X$-valued random variable

$$\mathcal{S}_{\mathcal{D}}(\mathsf{z}) + \triangle\mathsf{x} \quad \text{given a synthesis operator } \mathcal{S}_{\mathcal{D}} \colon Z \to X. \tag{6}$$

In the above, $\mathsf{z} \sim \mathsf{P}_Z$ is a $Z$-valued random variable generating dictionary coefficients and $\triangle\mathsf{x} \sim \mathsf{P}_{\triangle\mathsf{x}}$ is a $X$-valued random variable representing random variations that arise from limitations in the model capacity of the generative model given by the dictionary.

A well-chosen generative model needs to have sufficiently high model capacity in order to represent important features of the true signal, like edges and texture in an image at various scales. On the other hand, a too large model capacity will also represent noise and other unwanted features.

## 3.1   Dictionary learning as evidence maximisation

*Dictionary learning* is defined here as the task of choosing dictionary $\mathcal{D} \subset X$ so that the statistical distribution of generated signals is as close as possible to the true (unknown) distribution of the signal. More precisely, let $\mathsf{Q}_{\mathcal{D}}$ denote the distribution of signals generated as in eq. (6), i.e., $\mathcal{S}_{\mathcal{D}}(\mathsf{z}) + \triangle\mathsf{x} \sim \mathsf{Q}_{\mathcal{D}}$. We then seek the dictionary that maximises the match to the true distribution $\mathsf{P}_X$, e.g., by solving

$$\widehat{\mathcal{D}} \in \arg\min_{\mathcal{D}} \mathcal{W}\big(\mathsf{P}_X, \mathsf{Q}_{\mathcal{D}}\big). \tag{7}$$

Here, $\mathcal{W}$ quantifies similarity between probability distributions on $X$. The above essentially amounts to choosing the dictionary $\mathcal{D}$ so that $\mathcal{S}_{\mathcal{D}}(\mathsf{z}) + \triangle\mathsf{x}$ is as close as possible to $\mathsf{x}$ as random variables in $X$. The optimisation in eq. (7) can be non-convex, which is why we define $\widehat{\mathcal{D}}$ with '$\in$' instead of equality.

An issue with the formulation in eq. (7) is that it assumes one has access to the true (unknown) signal distribution $\mathsf{P}_X$. Often, one has access to i.i.d. samples $\Sigma := \{x_1, \ldots, x_N\} \subset X$ generated by $\mathsf{x} \sim \mathsf{P}_X$. One can replace the unknown distribution $\mathsf{P}_X$ in eq. (7) with $\widehat{\mathsf{P}}_{\Sigma} \in \mathscr{P}_X$, which is the empirical measure given by the aforementioned training data, so eq. (7) is replaced with

$$\widehat{\mathcal{D}} \in \arg\min_{\mathcal{D}} \mathcal{W}\big(\widehat{\mathsf{P}}_{\Sigma}, \mathsf{Q}_{\mathcal{D}}\big) \quad \text{where} \quad \widehat{\mathsf{P}}_{\Sigma} := \frac{1}{N} \sum_{i=1}^{N} \delta_{x_i}. \tag{8}$$

Here $\delta_{x_i}$ is the probability measure on $X$ that has a unit point mass at $x_i \in X$.

The next step is to specify $\mathcal{W} \colon \mathscr{P}_X \times \mathscr{P}_X \to \mathbb{R}$ in eq. (7) (and eq. (8)). The KL divergence is a common choice, so eq. (8) reads as

$$\widehat{\mathcal{D}} \in \arg\min_{\mathcal{D}} \mathrm{KL}\big(\widehat{\mathsf{P}}_{\Sigma} \mid \mathsf{Q}_{\mathcal{D}}\big). \tag{9}$$

If $\rho_{\mathcal{D}} \colon X \to \mathbb{R}_+$ is the density for $\mathsf{Q}_{\mathcal{D}}$, then the KL divergence is expressible as

$$\mathrm{KL}\big(\widehat{\mathsf{P}}_{\Sigma} \mid \mathsf{Q}_{\mathcal{D}}\big) := \frac{1}{N} \sum_{i=1}^{N} \Big\{ \log\Big(\frac{1}{N}\Big) - \log \rho_{\mathcal{D}}(x_i) \Big\}$$

$$= \log\Big(\frac{1}{N}\Big) - \frac{1}{N} \sum_{i=1}^{N} \log \rho_{\mathcal{D}}(x_i).$$

Hence, the dictionary learning problem in eq. (9) can be re-phrased as evidence maximisation:

$$\widehat{\mathcal{D}} \in \arg\max_{\mathcal{D}} \sum_{i=1}^{N} \log \rho_{\mathcal{D}}(x_i). \tag{10}$$

## 3.2   Sparse dictionary representations

To proceed, we introduce assumptions that are typical for sparse dictionary representations [45, 44, 32, 70].

**Assumption 3.1** (Sparse dictionary representation)**.** *Consider the generative model in eq.* (6) *where:*

*(A1) $X = \mathbb{R}^n$, $Z = \mathbb{R}^m$ with $n \ll m$.*

*(A2) Synthesis operator: $\mathcal{S}_{\mathcal{D}} \colon \mathbb{R}^m \to \mathbb{R}^n$ is linear, i.e.*

$$\mathcal{S}_{\mathcal{D}}(z) := \mathbf{D}z \quad where \ \mathbf{D} \in \mathbb{R}^{n \times m} . \tag{11}$$

*(A3) Dictionary atoms have a fixed norm:*

$$\|d\| = 1 \quad for \ all \ d \in \mathcal{D}. \tag{12}$$

*(A4) Dictionary coefficients are Laplace distributed, i.e., $z \sim \mathsf{P}_Z$ has corresponding density*

$$\rho_{\mathsf{z}}(z) = \frac{1}{(2b)^m} \exp\Big(-\frac{\|z\|_1}{b}\Big). \tag{13}$$

*(A5) $\triangle\mathsf{x}$ is independent of $\mathsf{z}$ and has Gaussian distribution with density*

$$\rho_{\triangle\mathsf{x}}(x) = \frac{1}{(\sqrt{2\pi}\sigma)^n} \exp\Big(-\frac{1}{2\sigma^2}\|x\|_2^2\Big). \tag{14}$$

Assumption (A1) merely states that the signal and latent spaces are both finite dimensional and the dictionary is over-complete, i.e., there are more dictionary atoms ($= m$) than the dimension of the signal ($= n$). Linearity assumption in (A2) implies that synthesis operator $\mathcal{S}_{\mathcal{D}}$ can be represented by a matrix. However, it does not necessarily mean that each dictionary atom in $\mathcal{D}$ is represented by a column vector of $\mathbf{D}$. As an example, we will use a convolutional synthesis operator where dictionary atoms (in $\mathcal{D}$) represent convolutional kernels and for each kernel there are columns in $\mathbf{D}$ that correspond to different shifts of that kernel. In addition, we fix the norm of dictionary atoms in assumption (A3). Without this assumption, one would have to take into account that distribution of the dictionary coefficients depends on the norm of corresponding dictionaries. This would make the model, in particular the assumption (A4), more complicated. The assumption (A4) ensures that maximizing the posterior density $z \mapsto \rho_{\mathsf{z}|\mathsf{x}}(z \mid x)$ yields sparsity of dictionary coefficients. Finally, in (A5) we use Gaussian distribution to account for small inaccuracies in the generative model $\mathcal{S}_{\mathcal{D}}(\mathsf{z})$.

We proceed by making use of the above assumptions for computing the objective (log-evidence) in eq. (10). These assumptions yield an expression

for the joint density for $(x, z)$, so the idea is to express the desired density by marginalizing over dictionary coefficients $z$:

$$\sum_{i=1}^{N} \log \rho_{\mathcal{D}}(x_i) = \sum_{i=1}^{N} \log \int_Z \rho_{\mathsf{x},\mathsf{z}}(x_i, z_i) \, \mathrm{d}z_i$$

$$= \sum_{i=1}^{N} \log \int_Z \rho_{\mathsf{x}|\mathsf{z}}(x_i \mid z_i) \rho_{\mathsf{z}}(z_i) \, \mathrm{d}z_i \qquad (15)$$

$$= \sum_{i=1}^{N} \log \int_Z \rho_{\triangle \mathsf{x}}(x_i - \mathcal{S}_{\mathcal{D}}(z_i)) \rho_{\mathsf{z}}(z_i) \, \mathrm{d}z_i.$$

Note that all density functions in the above expression except for $\rho_{\mathsf{z}}$ depend on dictionary $\mathcal{D}$, however we do not show this dependency to simplify the notation. Nevertheless, it is important to keep in mind that distributions $\rho_{\mathsf{x},\mathsf{z}}$ and $\rho_{\mathsf{x}|\mathsf{z}}$ are not the true distributions, but the ones assumed by the model (since we are expressing the log evidence $\log \rho_{\mathcal{D}}(x_i)$).

The $Z$-integrals above are computationally demanding. A standard machine learning approach to avoid this problem by decomposing the log evidence into two terms:

$$\log \rho_{\mathcal{D}}(x) = \mathrm{KL}(q \mid \rho_{\mathsf{z}|\mathsf{x}}) + \mathrm{ELBO}(q)(x) \geq \mathrm{ELBO}(q)(x) \qquad (16)$$

Note that $q \colon Z \to [0, 1]$ above is *any* probability density function on $Z$ and $\mathrm{ELBO}(q) \colon X \to \mathbb{R}$ is the evidence lower bound (ELBO) that is defined as

$$\mathrm{ELBO}(q)(x) = \int_Z q(z) \log \frac{\rho_{\mathsf{x},\mathsf{z}}(x, z)}{q(z)} \, \mathrm{d}z$$

$$= \mathbb{E}_q \big[ \log \rho_{\mathsf{x},\mathsf{z}}(x, \tilde{\mathsf{z}}) \big] - \mathbb{E}_q \big[ \log q(\tilde{\mathsf{z}}) \big], \qquad (17)$$

where $\tilde{\mathsf{z}}$ is any $Z$-valued random variable with a probability density function $q$ and $\mathbb{E}_q$ denotes the expectation w.r.t. this density. Next, the inequality in eq. (16) holds for any $q$ since the KL divergence is always positive. Hence, instead of directly maximizing the log evidence $\log \rho_{\mathcal{D}}(x)$ with respect to $\mathcal{D}$ as in eq. (10), one can maximize the ELBO. Furthermore, if $q$ approximates $\rho_{\mathsf{z}|\mathsf{x}}$, the KL divergence is small and therefore the lower bound is tight, see [13] for further details.

We consider $q$ being a Laplace distribution concentrated around the mode of $\rho_{\mathsf{z}|\mathsf{x}}$. This choice is supported by the following intuition: The mode of the posterior $z^*$ is a sparse vector. Therefore, most elements of its elements are 0. If the representation is accurate (the variance of $\triangle \mathsf{x}$ is small), then the dictionary model is capable of explaining a large part of the variability of $\mathsf{x}$ and $\rho_{\mathsf{x}|\mathsf{z}}(x \mid z)$ is larger than $\rho_{\mathsf{x}}(x)$ close to the mode $z^*$. If the representation is unique, the opposite holds further away from the mode. Therefore, along most of the dimensions the posterior

$$\rho_{\mathsf{z}|\mathsf{x}}(z \mid x) = \frac{\rho_{\mathsf{x}|\mathsf{z}}(x \mid z)}{\rho_{\mathsf{x}}(x)} \rho_{\mathsf{z}}(z) \qquad (18)$$

is even spikier than the prior $\rho_{\mathsf{z}}(z)$.

Furthermore, the following theorem shows that a slight relaxation of the lower bound will, given this approximation of the posterior $\rho_{\mathsf{z}|\mathsf{x}}$ by $q$, result in the optimization problem eq. (5). Typical approaches for dictionary learning are based on solving this optimisation, so the theorem below offers a statistical interpretation of the dictionary learning procedure. The theorem shows in particular that if the model is sparse, then solving eq. (5) essentially amounts to maximizing the ELBO for a Laplace distributed posterior. Moreover, if this posterior is appropriate (KL divergence with the true posterior is small), we are maximizing log-likelihood of image samples and hence minimizing the KL divergence between empirical distribution of images and distribution generated by our model.

**Theorem 3.1.** *Consider the generative model in eq.* (6) *and assume (A1)-(A5) holds. Also, let $q$ denote the density of a Laplace distributed $Z$-valued random variable centered at the mode for the posterior, i.e.*

$$q(z \mid x) = \frac{1}{(2b^*)^m} \exp\left(-\frac{\|z - z^*\|_1}{b^*}\right) \quad \text{for fixed } b^* > 0, \tag{19}$$

*with*

$$z^* := \arg\max_{z \in Z} \log \rho_{\mathsf{z}|\mathsf{x}}(x, z) = \arg\max_{z \in Z} \log \rho_{\mathsf{x},\mathsf{z}}(x, z). \tag{20}$$

*Then*

$$\begin{aligned} \mathrm{ELBO}(q)(x) &= -\mathbb{E}_q\big[f(x, \tilde{\mathsf{z}})\big] + C(\sigma, b, b^*) \\ &\geq -f(x, z^*) - m\frac{(b^*)^2}{\sigma^2} - m\frac{b^*}{b} + C(\sigma, b, b^*) \end{aligned} \tag{21}$$

*where*

$$\begin{aligned} f(x, z) &:= \frac{1}{2\sigma^2}\|\mathbf{D}z - x\|^2 + \frac{1}{b}\|z\|_1 \\ C(\sigma, b, b^*) &:= -\frac{n}{2}\log 2\pi\sigma^2 + m\log\frac{b^*}{b} + 1. \end{aligned} \tag{22}$$

*Finally, the gap between $\mathrm{ELBO}(q)$ and its lower bound in eq.* (21) *is bounded by $\frac{b^*}{b}\|z^*\|_0$.*

*Proof.* First, we express $\mathrm{ELBO}(q)$ as defined in eq. (17) by expressing its parts

$$\begin{aligned} \mathbb{E}_q\big[\log q(\tilde{\mathsf{z}} \mid x)\big] &= -m\log 2b^* - \frac{1}{b^*}\mathbb{E}_q\big[\|\tilde{\mathsf{z}} - z^*\|_1\big] \\ &= -m\log 2b^* - 1 \end{aligned} \tag{23}$$

and

$$\begin{aligned} \mathbb{E}_q\big[\log \rho_{\mathsf{x},\mathsf{z}}(x, \tilde{\mathsf{z}})\big] &= \mathbb{E}_q\Big[\log\big(\rho_{\mathsf{x}|\mathsf{z}}(x \mid \tilde{\mathsf{z}})\rho_{\mathsf{z}}(\tilde{\mathsf{z}})\big)\Big] \\ &= -\frac{n}{2}\log 2\pi\sigma^2 - m\log 2b - \mathbb{E}_q\big[f(x, \tilde{\mathsf{z}})\big] \end{aligned} \tag{24}$$

Set $\Delta\mathsf{z} := \tilde{\mathsf{z}} - z^*$, then

$$\begin{aligned} f(x, \tilde{\mathsf{z}}) &= \frac{1}{2\sigma^2}\big\|\mathbf{D}(z^* + \Delta\mathsf{z}) - x\big\|^2 + \frac{1}{b}\|z^* + \Delta\mathsf{z}\|_1 \\ &= \frac{1}{2\sigma^2}\|\mathbf{D}z^* - x\|^2 + \frac{1}{\sigma^2}(\mathbf{D}z^* - x)^\top \mathbf{D}\Delta\mathsf{z} \\ &\quad + \frac{1}{2\sigma^2}\|\mathbf{D}\Delta\mathsf{z}\|^2 + \frac{1}{b}\|z^* + \Delta\mathsf{z}\|_1 \end{aligned} \tag{25}$$

Now, we make the following observations:

- $\mathbb{E}_q[\Delta z] = \mathbb{E}_q[\tilde{z}] - z^* = 0$,

- $\mathbb{E}_q\big[\|\mathbf{D}\Delta z\|^2\big] = \mathbb{E}_q\big[\Delta z^\top \mathbf{D}^\top \mathbf{D} \Delta z\big] = 2(b^*)^2 \operatorname{tr}(\mathbf{D}^\top \mathbf{D}) = 2(b^*)^2 m$,

- $\|\tilde{z}\|_1$ is folded Laplace distribution and its expected value has been derived analytically in [46]:

$$\mathbb{E}_q\big[\|\tilde{z}\|_1\big] = \|z^*\|_1 + b^* \sum_{i=1}^m \exp\left(\frac{-|z^*_i|}{b^*}\right). \qquad (26)$$

Then, we can conclude that

$$
\begin{aligned}
\mathbb{E}_q\big[f(x, \tilde{z})\big] &= f(x, z^*) + m\frac{(b^*)^2}{\sigma^2} + \frac{b^*}{b}\sum_{i=1}^m \exp\left(\frac{-|z^*_i|}{b^*}\right) \\
&\leq f(x, z^*) + m\frac{(b^*)^2}{\sigma^2} + m\frac{b^*}{b}
\end{aligned}
\qquad (27)
$$

The desired lower bound on $\text{ELBO}(q)$ given in eq. (21) now follows from combining eq. (27) with eqs. (23) and (24). $\qquad\square$

We conclude with some remarks that relate to consequences of the theorem.

Note first that the gap $\frac{b^*}{b}\|z^*\|_0$ is relatively small. It constitutes a $\|z^*\|_0/m$ fraction of the term $mb^*/b$. Thus, with sufficient sparsity we are very close to optimizing ELBO even without sampling from the posterior.

Next, a similar result can be shown assuming that $q$ is a Gaussian density function, however in this case, the gap between ELBO and our optimization objective (lower bound of ELBO) might be larger. In any case, whether assumption on the posterior shape is appropriate or not is quantified by KL term.

Finally, this theorem also points to a connection between dictionary learning and variational autoencoders (VAs). The synthesis operator in dictionary learning can be seen as a decoder with one convolutional layer, the sparse coding defines an encoder. However, in dictionary learning the posterior is parametrized by assuming it is Laplace distributed around the mode, while VAs assume Gaussian distribution. Since Laplace distribution is very spiky a good (even though biased) approximation of ELBO is obtained without sampling from the posterior. Thus, the encoder does not have to predict variance of the posterior distribution, it is sufficient to predict the mode. In both cases the assumption that posterior has a certain form might not be appropriate, leaving the opportunity for both methods to fail.

## 3.3 Implementation of dictionary learning

In this section we discuss practical implementation of the synthesis operator $\mathcal{S}_\mathcal{D}\colon Z \to X$ and the procedure of learning a good dictionary $\mathcal{D}$.

### 3.3.1 Generative model

In the section 3 we made two assumptions regarding the generative model in eq. (6). First, the synthesis operator (A2) is linear, and second, dictionaries

have fixed norm (A3). In this section we present further choices that are made to apply the model in practice.

As commonly done, we use the generative model eq. (6) to model only high frequency features, while a low frequency component is estimated using FBP with a low-pass filter. Secondly, we use non-overlapping image patches to learn dictionary atoms. Further, follows the discussion of these choices.

**Subtracting the low frequency component**    The dictionary model eq. (6) needs to have sufficient capacity to represent "natural" signals in $X$. Digitizing such signals will in many applications, like imaging, result in high dimensional arrays. Sparsely representing the entire signal with a fixed dictionary will therefore require impractically many atoms. In addition, learning those atoms will be challenging, bearing in mind the limited number of training data examples that are available in practice. A common approach to address the above is to use dictionaries for representing only the high frequency component of the signal.

Using a dictionary of the above type in signal reconstruction must include a step that recovers the low frequency component of the signal from noisy data in order to remove it. For image de-noising, one can simply recover the low frequency component by low-pass filtering data (noisy image). The corresponding approach in CT image reconstruction is to use FBP with a cut-off frequency for the reconstruction kernel that is far below the Nyquist frequency that is dictated by sampling theory. This ensures an over-smoothed image representing the low frequency component. In our experiments, we set the cut-off frequency to 10% of the Nyquist frequency.

The above approach also requires one to train the dictionary against high frequency components of natural signals, i.e., one needs to remove low frequency components of images in the training set. This could be done by simply applying a low-pass filter to the images. However, in CT image reconstruction, we don't have images available and the low frequency component must be obtained differently. Using different procedures for extracting low frequency components during dictionary learning and during CT reconstruction is likely lead to sub-optimal results. For this reason, we employ the following procedure to remove low frequency components of images in the training set: First, generate noise free synthetic data by applying the CT forward operator on the high resolution CT images in the training set. Next, compute the corresponding low frequency component of the images by applying FBP on this noise free synthetic data. This results in an approach for dictionary learning adapted for CT reconstruction. Finally, the approach is unsupervised in the sense that it only requires access to high quality CT images, i.e., it does not require any access to CT projection data.

**Training on patches**    When we learn dictionary $\mathcal{D}$ we use a patch-based synthesis operator $\mathcal{S}_{\mathcal{D}}^{p} \colon Z \to X$ that generates an image from sparse representations of on non-overlapping image patches. In practice, we implement it as a convolutional operator with a stride that is equal to the size of a dictionary atom. However, when we use the learned dictionary for regularization, we set the stride to 1, i.e. use a regular convolutional synthesis operator $\mathcal{S}_{\mathcal{D}} \colon Z \to X$. We motivate this choice in a discussion below.

The generative model defined in eq. (1) traditionally has been applied to

image patches instead of whole images. To avoid "block" artifacts in signal restoration tasks (like de-noising and reconstruction), the whole image has been processed by applying the model to overlapping patches independently and then averaging the results.

Patch-based dictionary learning described above and convolutional dictionary learning define different ways to learn a dictionary. However, when a learned dictionary is applied for regularization in some task, the two approaches rely on very similar generative models. In a patch-based setting for each patch there is a corresponding sparse signal representation (dictionary coefficients). If these representations are combined together to form coefficient maps $z_i \in X$ for each dictionary element $d_i \in \mathcal{D}$, then a full image can be generated by applying a convolutional synthesis operator as defined eq. (2). In fact, the only difference between the two methods is that in the patch-based approach dictionary coefficients are computed independently for each patch. In contrast, in the convolutional framework the sparse coding is done jointly for all overlapping patches in the image.

Convolutional generative model have several advantages. Joint optimization with respect to coefficients allows to obtain closer approximation of an image given the same sparsity level, making the model more flexible. Moreover, a dictionary does not have to contain shifted versions of the same atoms, which makes the model more efficient [22]. On the other hand, in [65] authors argue that atoms in a convolutional dictionary suitable for representing natural images have high correlation with shifted versions of themselves, which makes uniqueness and stability guarantees for the sparse coding step obsolete. Non-uniqueness of the sparse representation is problematic in training/dictionary learning, because in this case the posterior distribution of coefficients is not unimodal and the Laplace approximation is not appropriate. Subsequently, solving dictionary learning problem eq. (5) might not maximize the log-likelihood eq. (16). This is supported by the observation that using convolutional dictionaries for regularization in de-noising has been problematic in practice. Indeed, works that we know have succeeded to learn only small (8) dictionary atoms for representing high frequency components.

Our initial goal was to use convolutional dictionary learning with large dictionary atoms, to provide stronger regularization in solving inverse problems. However, we also found that learning such atoms in convolutional setting is very hard, since it is hard to ensure both sparsity and that dictionary atoms are regularly updated. First, we tried to address this issue by dropout. However, we discovered that best option in terms of results and computational efficiency is to learn the dictionary on patches, by using a patch-bases synthesis operator $\mathcal{S}_{\mathcal{D}}^p$.

**Regularization with convolutional synthesis operator**   Since we discovered that it is beneficial to learn a dictionary on non-overlapping patches, it would be natural to apply this dictionary for regularization as it is usually done in the patch-based approach. However, we find that using the convolutional model with the learned dictionary gives slightly better results in practice, see appendix A for comparison. Therefore, we substitute operator $\mathcal{S}_{\mathcal{D}}^p$ by $\mathcal{S}_{\mathcal{D}}$. As a consequence, a different generative model is used in training and in application. A natural question is to elucidate how this affects the distribution $\mathsf{Q}_{\mathcal{D}}$ in section 3.1 that is given by the learned dictionary.

First, the new synthesis operator will increase the model capacity, therefore $x \mapsto f(x, z^*)$ in eq. (21) will decrease. Unfortunately, the model will also become better at representing unwanted features, which is reflected by an increase in the other terms in eq. (21), like the dimension $m$ of the coefficient space. This increase in model capacity most likely leads to overall decrease ELBO. However, the increase in model capacity is limited by the fact that during the training on non-overlapping patches, dictionaries are forced to learn shifted versions of themselves. Therefore, the posterior distribution of the coefficients becomes multi-modal and $\mathrm{KL}(q \mid \rho_{\mathsf{z}|\mathsf{x}})$ increases. This explains how the approach can still lead to good results in practice.

### 3.3.2   Solving the non-convex optimization problem

In principle one could learn the best generative model by maximizing the right-hand-side of eq. (21) jointly with respect to $\mathcal{D}$ and the hyper-parameters $\sigma$, $b$, and $b^*$. However, this results in a non-convex optimization problem whose solution depends on the initialization. More precisely, if one initializes with poor dictionaries and optimize with respect to hyper-parameters, one risks converging to a trivial solution where all image features are explained by noise and not the dictionary.

For the above reasons, we fix hyper-parameters and for given training data $\{x_i\}, i = 1, \ldots N$ (with low frequency components removed as outlined in section 3.3.1), we learn the dictionary by solving

$$\min_{\mathcal{D}} \sum_{i=1}^{N} \min_{z_i} \| \mathcal{S}_{\mathcal{D}}^{p}(z_i) - x_i \|^2 + \lambda \|z_i\|_1, \tag{28}$$

$$\|d\| = 1, \quad \text{for all } d \in \mathcal{D}.$$

**Stochastic optimization with alternating steps**   Since the number of available samples $N$ is high, we solve eq. (28) using stochastic optimization with batch size 1. This means that at each iterate $k$, only one random image $x_{i_k}$ is used to update the dictionary. In addition, we randomly crop a smaller image patch from the image.

We solve the optimization problem by minimizing the objective with respect to dictionary and coefficients in alternating steps. First, dictionary coefficients $z_k$ are computed by solving a convex sub-problem:

$$z_k = \arg\min_{z} \big\| \mathcal{S}_{\mathcal{D}_k}(z) - x_{i_k} \big\|_2^2 + \lambda_k \|z\|_1. \tag{29}$$

This is done by applying fast ISTA (FISTA) [11] for a fixed number of iterates. Secondly, the dictionary is updated using the previous iterate $\mathcal{D}_k$ and the $\mathcal{D}$-gradient of the objective in eq. (29) at $\mathcal{D}_k$:

$$\mathcal{D}_{k+1} = \mathrm{OptUpdate}\bigg( \mathcal{D}_k, \nabla_{\mathcal{D}} \| \mathcal{S}_{\mathcal{D}_k}(z_k) - x_{i_k} \|_2^2 \bigg). \tag{30}$$

In particular, simple stochastic gradient descent (SGD) has updates of the form

$$\mathcal{D}_{k+1} = \mathcal{D}_k - \alpha_k \nabla_{\mathcal{D}} \| \mathcal{S}_{\mathcal{D}_k}(z_k) - x_{i_k} \|_2^2. \tag{31}$$

We find that using the Adam optimizer [40] instead of SGD leads to better results. In Adam the accumulated gradient is normalized with respect to accumulated variance of the gradients. In dictionary learning, this helps to avoid so called "dead" atoms - the atoms that are not picked up in the sparse coding step and subsequently not updated. Indeed, we learn two times more different dictionary atoms while using the Adam optimizer compared to regular SGD.

Note that $\nabla_{\mathcal{D}}$ in eqs. (30) and (31) is the $\mathcal{D}$-gradient of $\mathcal{D} \mapsto \left\| \mathcal{S}_{\mathcal{D}}(z) - x \right\|_2^2$ which can be explicitly expressed as

$$\nabla_{\mathcal{D}} \left\| \mathcal{S}_{\mathcal{D}}(z) - x \right\|_2^2 = 2 \left[ \partial_{\mathcal{D}} \, \mathcal{S}_{\mathcal{D}}(z) \right]^* \left( \mathcal{S}_{\mathcal{D}}(z) - x \right).$$

Here, $\partial_{\mathcal{D}} \, \mathcal{S}_{\mathcal{D}}(z)$ is the $\mathcal{D}$-Jacobian of $\mathcal{D} \mapsto \mathcal{S}_{\mathcal{D}}(z)$.

**Adaptive choice of regularization parameter**   If the regularization parameter $\lambda$ is constant during the training, many dictionary atoms will not be updated in the beginning of the optimization process and therefore will not be used later. Thus, it is important to be flexible with sparsity level $\lambda$ to insure that majority of the atoms get a good start. We address this problem by adopting $\lambda$ during the optimization so that the average number of non-zero elements in $z$ stays approximately constant. This results in $\lambda$ rising in the beginning of the optimization and then stabilizing around one value.

This is done by estimating the average sparsity level on the validation set $\widehat{s} = \frac{1}{N_{val}} \sum_{i=1}^{N_{val}} \| z_i^t \|_0$. If $s$ is the pre-set sparsity level, then during every validation step we have for some small constant $c$

$$\lambda_{t+1} = \begin{cases} \lambda_t + c(\widehat{s} - s) & \text{if } |\widehat{s} - s| > 0.2s, \text{ and } t \bmod 10 = 0 \\ \lambda_t & \text{otherwise.} \end{cases} \tag{32}$$

An alternative approach to dictionary initialization is to apply clustering algorithms to image patches and initialize atoms as cluster centres [3, 7]. However, this approach is likely to be too demanding from the computational point of view, given a large number of image patches that we have.

## 4   Dictionary based regularization in inverse problems

An inverse problems is a problem of recovering an unknown signal $x^* \in X$ from data $y \in Y$ that represents indirect noisy observations of the signal. In statistical terms we are interested in the distribution of the conditional random variable $(\mathsf{x} \mid \mathsf{y} = y)$, where

$$\mathsf{y} = \mathcal{A}(\mathsf{x}) + \mathsf{e}. \tag{33}$$

Here, the mapping $\mathcal{A} \colon X \to Y$ (forward operator) is assumed to be known and $\mathsf{e}$ is a random variable that generates observation errors.

Given a prior distribution $\mathsf{x} \sim \mathsf{P}_X$ the posterior density is expressible as

$$\rho_{\mathsf{x}|\mathsf{y}}(x \mid y) \propto \rho_{\mathsf{y}|\mathsf{x}}(y \mid x)\rho_{\mathsf{x}}(x) = \rho_{\mathsf{e}|\mathsf{x}}\big(y - \mathcal{A}(x) \mid x\big)\rho_{\mathsf{x}}(x) \tag{34}$$

where the density $\rho_{\mathsf{e}|\mathsf{x}}$ for the observation error is usually know. Furthermore, the prior $\mathsf{P}_X$ can be approximated by $\mathsf{Q}_{\widehat{\mathcal{D}}}$, where $\widehat{\mathcal{D}}$ is a learned dictionary:

$$\rho_{\mathsf{x}}(x) \approx \int \rho_{\mathsf{x}|\mathsf{z}}(x \mid z)\rho_{\mathsf{z}}(z) \, \mathrm{d}z. \tag{35}$$

The true solution $x^*$ to the inverse problem eq. (33) can be estimated by the maximum a posteriori (MAP) estimator, which is defined as the maximum of $x \mapsto \rho_{\mathsf{x}|\mathsf{y}}(x \mid y)$. Assuming the true prior can be approximated by a learned dictionary as in eq. (35) introduced a marginalization over $\mathsf{z}$, which is computationally unfeasible, so an alternative is to compute unmarginalized MAP estimator:

$$
\begin{aligned}
(\widehat{x}, \widehat{z}) \in \;& \underset{(x,z)\in X\times Z}{\arg\max} \; \rho_{\mathsf{x},\mathsf{z}|\mathsf{y}}(x, z \mid y) \\
=\;& \underset{(x,z)\in X\times Z}{\arg\max} \; \rho_{\mathsf{y}|\mathsf{x}}(y \mid x)\rho_{\mathsf{x}|\mathsf{z}}(x \mid z)\rho_{\mathsf{z}}(z) \\
=\;& \underset{(x,z)\in X\times Z}{\arg\max} \; \rho_{\mathsf{e}|\mathsf{x}}\big(y - \mathcal{A}(x) \mid x\big)\rho_{\triangle\mathsf{x}}(x - \mathcal{S}_{\widehat{\mathcal{D}}}(z))\rho_{\mathsf{z}}(z)\,\mathrm{d}z.
\end{aligned}
\tag{36}
$$

Inserting the specific expressions for the densities of the random variables $(\mathsf{e} \mid \mathsf{x})$, $\triangle\mathsf{x}$, and $\mathsf{z}$ into eq. (36) yields

$$
(\widehat{x}, \widehat{z}) \in \underset{(x,z)\in X\times Z}{\arg\max} \; \mathcal{L}\big(\mathcal{A}(x), y\big) + \lambda_1\big\|x - \mathcal{S}_{\widehat{\mathcal{D}}}(z)\big\|_2^2 + \lambda_2\|z\|_1
\tag{37}
$$

where $\mathcal{L}\big(\mathcal{A}(x), y\big) := -\log \rho_{\mathsf{e}|\mathsf{x}}\big(y - \mathcal{A}(x) \mid x\big)$ is the data (negative) log-likelihood and the regularization parameters are $\lambda_1 := 1/\sigma$ and $\lambda_2 := 1/b$. An equivalent bi-level optimization formulation to eq. (37) reads as

$$
\begin{cases}
\widehat{x} \in \underset{x\in X}{\arg\min} \; \mathcal{L}\big(\mathcal{A}(x), y\big) + \lambda_1\, \mathcal{R}_{\widehat{\mathcal{D}}}(x) & (38) \\[2mm]
\mathcal{R}_{\mathcal{D}}(x) \in \underset{z\in Z}{\min}\big\|x - \mathcal{S}_{\mathcal{D}}(z)\big\|_2^2 + \dfrac{\lambda_2}{\lambda_1}\|z\|_1. & (39)
\end{cases}
$$

## 4.1 Implementation

Signals (images) and measured data are in practice digitized and represented by arrays, i.e., $X = \mathbb{R}^n$ and $Y = \mathbb{R}^l$ in eq. (33). Next, we assume that observation noise $(\mathsf{e} \mid \mathsf{x} = x)$ has a Poisson distribution conditioned on noise-free data $\mathcal{A}(x^*)$ generated by the signal. To simplify computations we use a quadratic approximation for the data log-likelihood in eq. (37), that results in the following weighted $l_2-$norm [27]:

$$
\mathcal{L}\big(\mathcal{A}(x), y\big) = \sum_{i=1}^{l} w_i\big\|\mathcal{A}(x)_i - y_i\big\|_2^2 \quad \text{where weights } w_i = e^{-y_i}.
\tag{40}
$$

There are now several ways to solve the corresponding convex optimisation problem in eq. (37). First, FISTA [11] can be formulated with $(x, z)$ as the control variable. A different approach is to use an alternating scheme and apply a gradient descent step to variable $x$ and a proximal gradient descent to $z$ in alternating steps. Such a scheme would guarantee a monotonic decrease of the objective function given an appropriate selection of the step size [12]. In order to speed up the optimization, we use an accelerated gradient descent [54] to update $x$, and an accelerated proximal gradient descent (as in FISTA [11]) to update $z$,

which results in the following iterative scheme:

$$t_{k+1} = \frac{1 + \sqrt{1 + 4t_k^2}}{2}$$

$$x_{k+1} = x_k' - \frac{1}{L_x}\nabla_x f(x_k', z_k')$$

$$x_{k+1}' = x_{k+1} + \frac{t_k - 1}{t_{k+1}}\left(x_{k+1} - x_k\right) \tag{41}$$

$$z_{k+1} = \mathrm{prox}_g^{1/L_z}\left(z_{k+1}' - \frac{1}{L_z}\nabla_z f(x_{k+1}, z_k')\right)$$

$$z_{k+1}' = z_{k+1} + \frac{t_k - 1}{t_{k+1}}\left(z_{k+1} - z_k\right)$$

for $k \geq 1$ where $f(x, z) := \mathcal{L}\big(\mathcal{A}(x), y\big) + \lambda_1\big\|x - \mathcal{S}_{\widehat{\mathcal{D}}}(z)\big\|_2^2$ is the smooth part of the objective in eq. (37) and $L_x, L_z$ are the Lipschitz constants of the gradient of $f$ with respect to $x$ and $z$. The proximal operator for the non-smooth part of the objective $g(z) = \lambda_2\|z\|_1$ admits a closed form expression

$$\mathrm{prox}_g^\gamma(u) = \arg\min g(z) + \frac{1}{2\gamma}\|z - u\|^2$$

$$= \mathrm{sign}(u)\max\big(|u| - \gamma\lambda_2, 0\big).$$

The essential difference between the above scheme and FISTA (for the variable $(x, z)$) is that we are separately estimating the step size for $x$ and $z$, which results in using larger step sizes and, hence, faster convergence. However, we are not aware of any theoretical results guaranteeing convergence of our accelerated scheme in the general case, we verify that it results in a monotonic decrease of the objective function values in our experiments. Note, that we perform the optimization for the limited number of iterations, thus we do not necessarily reach the state of full convergence. In fact, running the process until full convergence leads to sub-optimal results. This observation that early stopping provides additional regularisation has been seen with other variational methods, like TV regularisation [26].

## 5   Application to CT image reconstruction

### 5.1   Experimental setting

To perform the experiments we simulate projection data using images of human abdomen provided by the 2016 AAPM Low Dose CT Grand Challenge [49]. The dataset contains CT scans of 10 patients, 9 of which are used for training and one is reserved for testing. To avoid computational burden we work in a 2D setting, splitting 3D volumes into axial slices of size $512 \times 512$. This results in 2168 and 210 image-data pairs for training and testing respectively. Moreover, 1% of the training images are separated into validation set, which is used for setting the hyper-parameters for the evaluated methods. The choice of hyper-parameters is made with an objective to minimize the mean squared error (MSE) on the validation set. Lastly, we compute evaluation metrics as for one 3D image, when we report the results, even though the methods are applied in a slice by slice manner.

Since the images are given in Hounsfield scale, we re-scale those by a factor $\mu_0/1000$, where $\mu_0 = 0.0192\text{mm}^{-1}$ is X-ray attenuation of water at the mean X-ray energy of 70 keV. The noise-free tomographic data $y = \mathcal{A}(x)$ is simulated using cone beam ray transform with 1000 detector elements and 1000 projection angles. According to the Beer–Lambert law the number of photons absorbed by the detector is

$$N_d = N_0 e^{-\mathcal{A}(x)}. \tag{42}$$

Here $N_0$ is an average number of photons that would reach a detector element in an empty space. Since $N_0$ is proportional to radiation intensity, smaller values correspond to a lower dose. For a single energy X-ray noise is modelled as a Poisson distributed random variable with mean $N_d$:

$$N_{\text{noisy}} = \text{Poisson}\left(N_0 e^{-\mathcal{A}(x)}\right) \tag{43}$$

Finally, we linearize the noisy data

$$y_{\text{noisy}} = -\ln\left(\frac{1}{N_0} N_{\text{noisy}}\right) \tag{44}$$

so that $\mathcal{A}(x) \approx y_{\text{noisy}}$. In our experiments we use $N_0 = 50\,000$, which corresponds to a noise level experienced in clinical low-dose CT imaging.

We learn 512 dictionary atoms of size $16 \times 16$ using optimization procedure described in section 3.3. During the training we aim for average sparsity level $\frac{3}{512} \approx 0.6\%$. During the reconstruction we set $\lambda_1 = 50$ and $\lambda_2 = 0.0016$ obtained by minimizing the validation error. The learned dictionary is visualized in appendix B.

First of all, we compare our method against FBP, which is an analytic reconstruction method that computes a regularised approximate inverse. We use the (approximate) implementation of FBP in Operator Discretisation Library (ODL) [1] with the Hanning filter and a relative frequency cut-off 0.75.

**Variational models**   Our comparison includes reconstruction by several variational models. All models can be defined as a solution to

$$\widehat{x} = \arg\min_{x \in X} \mathcal{L}\big(\mathcal{A}(x), y\big) + \lambda\, \mathcal{R}(x) \tag{45}$$

with $\mathcal{L}\colon Y \times Y \to \mathbb{R}_+$ as in eq. (40) and with different choices of $\mathcal{R}\colon X \to \mathbb{R}$.

The hyper-parameter $\lambda$ is a regularization parameter that governs the trade-off between stability and the need to fit data. Ideally the reconstruction methods includes a parameter selection rule for setting it, typically based on the noise level in data. These do not necessarily ensure best performance, so in order to ensure the different methods are compared at their best performance, we set this parameter empirically against validation data as to optimise performance. This also applies to the number of iterations used in the optimization scheme for solving eq. (45). A discussion on the potential regularising property of early stopping for variational methods is given in [26].

One variational model is TV [62], which corresponds to choosing $\mathcal{R}(x) = |\nabla x|_1$ in eq. (45) with $\nabla$ denoting the spatial gradient operator. The optimal choice of the regularisation parameter is $\lambda = 3 \cdot 10^{-4}$. We then solve eq. (45) using

the alternating direction method of multipliers (ADMM) algorithm with $1\,000$ iterations.

The next variational model is total generalised variation (TGV), which was introduced partly to address some of the drawbacks (like stair-casing) that comes with using TV [17, 16]. TGV corresponds to selecting

$$\mathcal{R}(x) = \min_{z} ||\nabla x - z||_1 + \alpha ||\mathcal{E} z||_1 \tag{46}$$

in eq. (45) where $\mathcal{E}$ is the symmetrized gradient operator, see [16] for the details. The optimal choice of regularisation parameters in our case is $\lambda = 6 \cdot 10^{-4}$ and $\alpha = 0.5$. The same optimization procedure as for TV is then used to solve eq. (45) with eq. (46).

Finally, we also include the commonly used Huber regularization, which is a smooth relaxation of TV regularization that replaces the non-smooth $l_1$-norm with the smooth Huber functional [28]. This corresponds to $\mathcal{R}(x) = \mathcal{H}_\gamma(\nabla x)$ in eq. (45) with $\mathcal{H}_\gamma \colon \mathbb{R}^{2n} \to \mathbb{R}_+$ (Huber functional) defined as

$$\mathcal{H}_\gamma(z) = \sum_{i=1}^{2n} \frac{z_i^2}{2\gamma} \mathbf{1}_{|z_i| < \gamma} + \sum_{i=1}^{2n} \left( |z_i| - \frac{\gamma}{2} \right) \mathbf{1}_{|z_i| \geq \gamma} \tag{47}$$

The optimal choice of regularisation parameters for our case is $\lambda = 5 \cdot 10^{-4}$ and $\gamma = 4 \cdot 10^{-4}$. The optimization eq. (45) with eq. (47) is solved using Nesterov's accelerated gradient descent [54] with 70 iterations.

**Deep learning based methods**   Deep image prior (DIP) relies on the fact that merely a representation of a signal as an output of a convolutional neural network $x = f_\theta(z)$ is sufficient for regularization in inverse problems [72, 8]. Here, we use an architecture of a neural network suggested by [72] for denoising to perform image reconstruction by minimizing

$$\min_{\theta} \mathcal{L}\big(\mathcal{A}(f_\theta(z)), y\big) \tag{48}$$

with respect to parameters of the neural network. The latent variable $z = z_0 + \Delta z$, where $z_0$ is sampled only once (before the start of the optimization process) and $\Delta z$ is regularizing noise, which is sampled during every iteration for solving eq. (48). This approach is completely unsupervised, since it does not require any training prior to the reconstruction. We use the same optimization procedure and hyper-parameters, with few exceptions. We reduce the standard deviation of the regularizing noise by a factor of 2 and increase the number of iterations to 6000. We observe that the approach might slightly benefit from additional iterations an higher noise, however even in this setting it is already remarkably slow compared to the other methods.

Another learned method that we used for comparison is adversarial regularization (AR) proposed in [47]. The method uses variational formulation eq. (45), where the regularizing component $\mathcal{R}(x)$ is a trained neural network. This method is can be seen as a semi-supervised, since it is trained using samples of both, reconstructions and projection data. However it does not rely on samples being coupled, i.e. correspond to the same scan. We use a default convolutional architecture and set regularization parameter $\lambda$ as suggested by the authors. The optimizations is done for 2000 iterations.

Finally, learned primal-dual (LPD) is a supervised deep learning method that approximately computes the posterior mean image given measured data, so it has no regularization parameter. It has shown state-of-the-art performance [2] in low-dose CT reconstruction and we include a version of this method applied to log-data.

We have paid a reasonable amount of effort to improve the performance of all schemes in the particular test setting. In particular, we re-scaled the images for the deep learning methods so that the gray-scale values approximately fit the interval $[0, 1]$. This is important, because the initialization of the neural networks is traditionally adapted for this scaling. Unsurprisingly, the architectures and hyper-parameters proposed by the authors in most cases provided the optimal performance. This observation is supported by the fact that [47] and [2] has been originally evaluated on the data simulated from same image dataset.

## 5.2   Results

Quantitative performance results in terms of peak signal to noise ratio (PSNR) and structural similarity index (SSIM) are given in table 1. In both metrics, LPD performs best. Second best is dictionary learning (DL) followed by AR. FBP performs worst in both metrics among the tested algorithms.

Figures 1 and 2 show exemplary slices from abdomen and pelvis, respectively, for a qualitative assessment of the image quality.

Figure 1 highlights two region of interests, which are well suited to assess the noise appearance and the delineation of low-contrast structures.

With respect to noise texture, TV shows the well-known behaviour with unnaturally looking flat regions with remaining salt-and-pepper noise. TGV shows a horizontally structured noise pattern, which is also present (but less pronounced) in the image with Huber regularization. AR shows a rather natural noise pattern, except for the aspect that it appears more low-frequent than the noise pattern in ground truth and FBP images. The most natural noise pattern is achieved by DIP and LPD. Finally, in DL barely any noise is remaining.

Regarding the low-contrast delineation, we focus on the boundaries of the pancreas in the lower right zoomed image region. Here, TV shows the well-known scruffy edges. The edge appearance improves in TGV at the cost of a higher noise level. Edges appear visibly blurred in DIP and AR, which is consistent for AR with the observation of the presence of rather low-frequent noise. Pancreas delineation is best for Huber, DL and LPD.

Figure 2 highlights the performance of the algorithms on low-contrast structures in the fatty tissue (bottom left zoom) high-contrast structures (bottom right zoom).

With respect to the structures in the fatty tissue, TV appears patchy and unnaturally looking. All other methods show a comparable level of details. More differences between the algorithms are present in the bony structure inside the hip joint shown in the zoom in the bottom right. Note the arc-shaped bony structure in the center of the ground truths (fig. 2a). This structure is not recovered at all in DIP and AR, while it is only partially visible in AR, Huber, TV, TGV. Only LPD and DL properly recover this structure.

In summary, the qualitative image quality assessment supports the quantitative evaluation that LPD performs best followed by DL.

(a) True        (b) FBP        (c) TV

(d) TGV        (e) Huber        (f) DIP
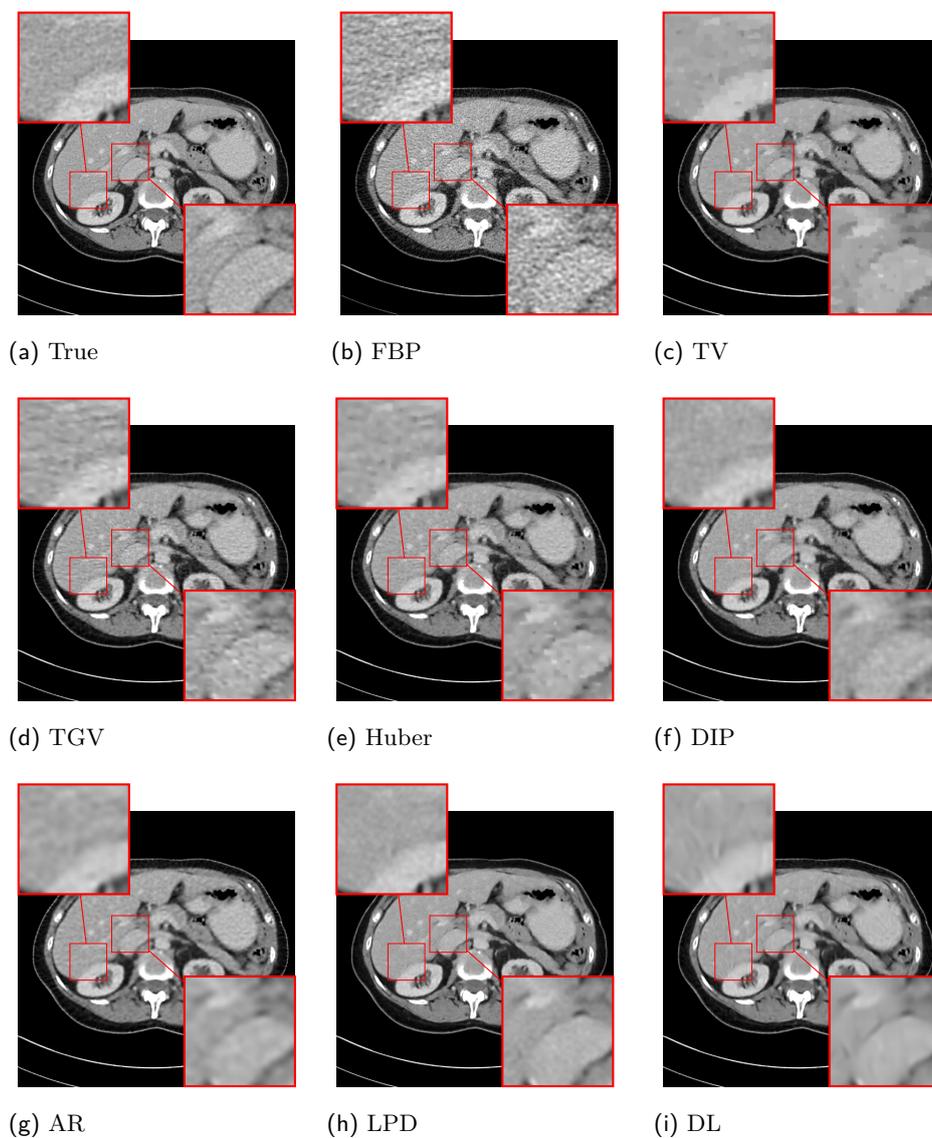
(g) AR        (h) LPD        (i) DL

Fig. 1: Example slice of the abdomen (level 80 HU, window 370 HU). Region of interests illustrate noise texture (top left) and delineation of low contrast structures (liver-kidney border in the top left region of interest and liver-pancreas transition as well as contrasted vessel in the liver in the bottom right region of interest.

|      | FBP   | TV    | TGV   | Huber | DIP   | AR    | LPD       | DL    |
|------|-------|-------|-------|-------|-------|-------|-----------|-------|
| PSNR | 40.95 | 46.37 | 46.70 | 46.70 | 46.29 | 46.96 | **49.68** | 48.20 |
| SSIM | 0.942 | 0.987 | 0.987 | 0.988 | 0.988 | 0.988 | **0.992** | 0.990 |

Tab. 1: Performance metrics for various reconstruction methods in low-dose CT. Note that LPD requires supervised training data, whereas DL can be trained against unsupervised data.

(a) True

(b) FBP

(c) TV

(d) TGV

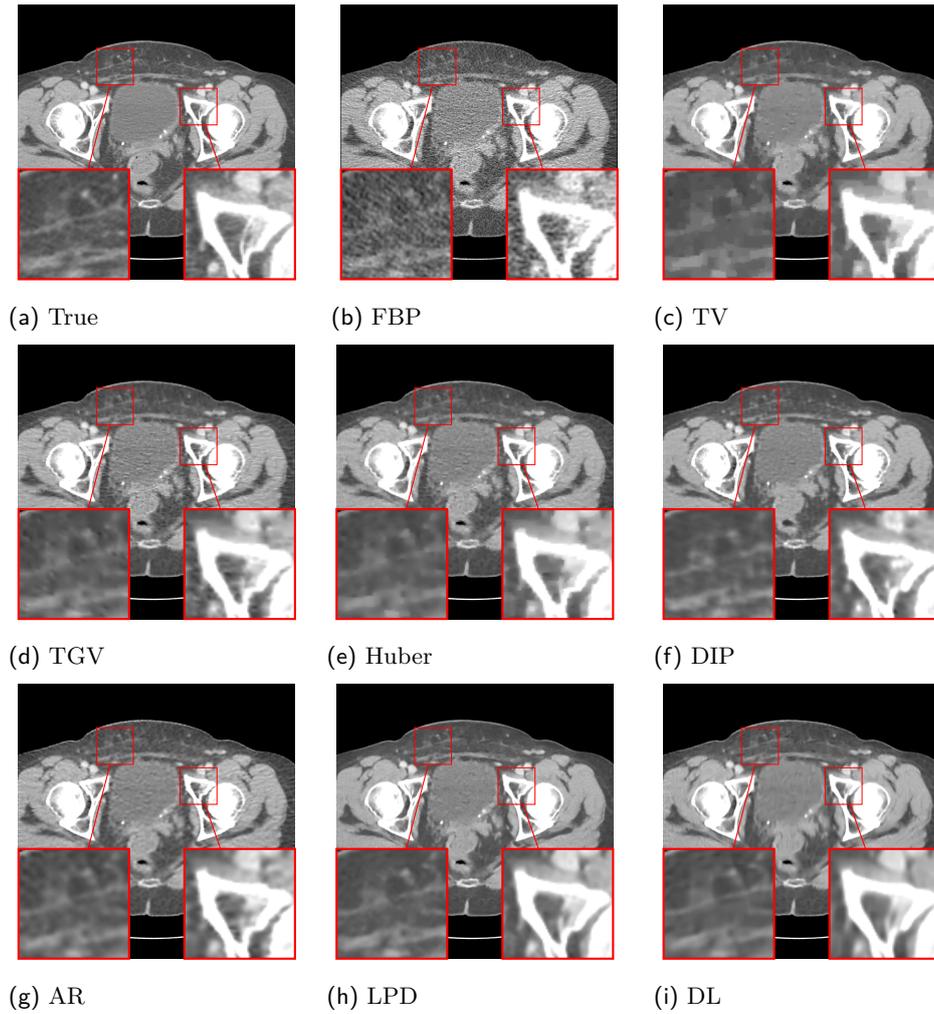(e) Huber

(f) DIP

(g) AR

(h) LPD

(i) DL

Fig. 2: Example slice of the pelvis (level 0 HU, window 400 HU). Regions of interest highlight the structure in the fatty tissue (bottom left) and fine high-contrast structures inside the hip joint (bottom right).

## 6   Conclusion

In this work we have presented a novel view point on the dictionary learning, by showing that it maximizes the evidence lower bound similarly to variational auto-encoders. Moreover, we have shown that dictionaries can be successfully learned using optimization techniques for training neural networks. We verify that regularization learned dictionaries constitutes a powerful method for tomographic reconstruction, which is unsupervised with respect to tomographic data. We compare our method with many different model-based and data-driven approaches and conclude that it outperforms most of the methods by a large margin, while being inferior only to the learned primal-dual, which is a fully supervised deep-learning method.

Even though dictionary learning based reconstruction leads to high quality images in terms of PSNR, we observe that this method suffers from high frequency noise patterns shaped by dictionary atoms. On the other hand, the methods that rely on multi-layer convolutional neural networks, such as LPD and DIP, produce noise patterns with lower frequency. Therefore, we suppose that a more complex generative model should be able to provide further improvements in image quality.

Finally, all evaluated variational methods, including our dictionary learning based approach, required many iterations to converge. This compromises their applicability in practice. Addressing this problem by learned optimization [9] is another direction for future research.

## Acknowledgments

## References

[1]   J. Adler, H. Kohr, and O. Öktem. *Operator Discretization Library (ODL)*. 2017. URL: https://github.com/odlgroup/odl/.

[2]   J. Adler and O. Öktem. "Learned primal-dual reconstruction". In: *IEEE transactions on medical imaging* 37.6 (2018), pp. 1322–1332.

[3]   A. Agarwal, A. Anandkumar, and P. Netrapalli. "Exact recovery of sparsely used overcomplete dictionaries". In: *stat* 1050 (2013), pp. 8–39.

[4]   M. Aharon, M. Elad, and A. Bruckstein. "K-SVD: An algorithm for designing overcomplete dictionaries for sparse representation". In: *IEEE Transactions on signal processing* 54.11 (2006), pp. 4311–4322.

[5]   E. Amaldi and V. Kann. "On the approximability of minimizing nonzero variables or unsatisfied relations in linear systems". In: *Theoretical Computer Science* 209.1 (1998), pp. 237–260.

[6]   H. K. Andersen, D. Völgyes, and A. C. T. Martinsen. "Image quality with iterative reconstruction techniques in CT of the lungs — A phantom study". In: *European Journal of Radiology Open* 5 (2018), pp. 35–40.

[7]    S. Arora, R. Ge, and A. Moitra. "New algorithms for learning incoherent and overcomplete dictionaries". In: *Conference on Learning Theory*. PMLR. 2014, pp. 779–806.

[8]    D. O. Baguer, J. Leuschner, and M. Schmidt. "Computed Tomography Reconstruction Using Deep Image Prior and Learned Reconstruction Methods". In: *arXiv preprint* 2003.04989 (2020).

[9]    S. Banert et al. "Accelerated Forward-Backward Optimization using Deep Learning". In: *arXiv preprint* (2021).

[10]   P. Bao et al. "Convolutional sparse coding for compressed sensing CT reconstruction". In: *IEEE transactions on medical imaging* 38.11 (2019), pp. 2607–2619.

[11]   A. Beck and M. Teboulle. "A Fast Iterative Shrinkage-Thresholding Algorithm for Linear Inverse Problems". In: *SIAM Journal on Imaging Sciences* 2.1 (2009), pp. 183–202.

[12]   A. Beck. *First-order methods in optimization*. SIAM, 2017.

[13]   C. M. Bishop. *Pattern recognition and machine learning*. springer, 2006.

[14]   T. Blumensath and M. E. Davies. "Iterative thresholding for sparse approximations". In: *Journal of Fourier analysis and Applications* 14.5-6 (2008), pp. 629–654.

[15]   R. N. Bracewell and A. C. Riddle. "Inversion of fan-beam scans in radio astronomy". In: *Astrophysical Journal* 150 (1967), pp. 427–434.

[16]   K. Bredies and M. Holler. "A TGV-based framework for variational image decompression, zooming, and reconstruction. Part II: Numerics". In: *SIAM Journal on Imaging Sciences* 8.4 (2015), pp. 2851–2886.

[17]   K. Bredies, K. Kunisch, and T. Pock. "Total generalized variation". In: *SIAM Journal on Imaging Sciences* 3.3 (2010), pp. 492–526.

[18]   D. J. Brenner and E. J. Hall. "Computed Tomography – An Increasing Source of Radiation Exposure". In: *New England Journal of Medicine* 357 (2007), pp. 2277–2284.

[19]   E. Candes, J. Romberg, and T. Tao. "Stable signal recovery from incomplete and inaccurate measurements". In: *Communications on pure and applied mathematics* 59.9 (2006), pp. 1207–1223.

[20]   E. Candes and T. Tao. "Near-optimal signal re- covery from random projections: Universal encoding strategies?" In: *IEEE Transaction on Information Theory* 52.12 (2006), pp. 5406–5425.

[21]   E. J. Candès, J. Romberg, and T. Tao. "Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information". In: *IEEE Transactions on information theory* 52.2 (2006), pp. 489–509.

[22]   A. Chambolle, M. Holler, and T. Pock. "A convex variational model for learning convolutional image atoms from incomplete data". In: *Journal of Mathematical Imaging and Vision* 62.3 (2020), pp. 417–444.

[23]   H. Chen et al. "LEARN: Learned experts' assessment-based reconstruction network for sparse-data CT". In: *IEEE transactions on medical imaging* 37.6 (2018), pp. 1333–1347.

[24]    Y. Chen et al. "Artifact suppressed dictionary learning for low-dose CT image processing". In: *IEEE transactions on medical imaging* 33.12 (2014), pp. 2271–2292.

[25]    D. Donoho. "For most large underdetermined systems of linear equations the minimal $\ell_l$-norm solution is also the sparsest solution". In: *Communications on pure and applied mathematics* 59.6 (2006), pp. 797–829.

[26]    A. Effland et al. "Variational Networks: An Optimal Control Approach to Early Stopping Variational Methods for Image Restoration". In: *Journal of Mathematical Imaging and Vision* 62 (2020), pp. 396–416.

[27]    I. A. Elbakri and J. A. Fessler. "Statistical image reconstruction for polyenergetic X-ray computed tomography". In: *IEEE transactions on medical imaging* 21.2 (2002), pp. 89–99.

[28]    H. Erdoğan. *Statistical image reconstruction algorithms using paraboloidal surrogates for PET transmission scans.* University of Michigan, 1999.

[29]    V. Faber, A. I. Katsevich, and A. G. Ramm. "Inversion of cone-beam data and helical tomography". In: *Journal of Inverse and Ill-posed Problems* 3 (Jan. 1995), pp. 429–446.

[30]    L. A. Feldkamp, L. C. Davis, and J. W. Kress. "Practical cone-beam algorithm". In: *Journal of the Optical Society of America A* 1.6 (1984), pp. 612–619.

[31]    L. L. Geyer et al. "State of the Art: Iterative CT Reconstruction Techniques". In: *Radiology* 276.2 (2015), pp. 339–357.

[32]    M. Girolami. "A variational method for learning sparse and overcomplete representations". In: *Neural computation* 13.11 (2001), pp. 2517–2532.

[33]    Z. Hu et al. "Image reconstruction from few-view CT data by gradient-domain dictionary learning". In: *Journal of X-ray science and technology* 24.4 (2016), pp. 627–638.

[34]    K. H. Jin et al. "Deep convolutional neural network for inverse problems in imaging". In: *IEEE Transactions on Image Processing* 26.9 (2017), pp. 4509–4522.

[35]    A. Kak and M. Slaney. *Principles of Computerized Tomographic Imaging.* Vol. 33. Classics in Applied Mathematics. Society for Industrial and Applied Mathematics, 2001. DOI: 10.1137/1.9780898719277.

[36]    M. K. Kalra et al. "Strategies for CT Radiation Dose Optimization". In: *Radiology* 230.3 (2004).

[37]    H. Kamoshita et al. "Low-Dose CT Reconstruction with Multiclass Orthogonal Dictionaries". In: *2019 IEEE International Conference on Image Processing (ICIP)*. IEEE. 2019, pp. 2055–2059.

[38]    A. Katsevich. "3PI algorithms for helical computer tomography". In: *Advances in Applied Mathematics* 36.3 (2006), pp. 213–250.

[39]    A. Katsevich. "A general scheme for constructing inversion algorithms for cone beam CT". In: *International Journal of Mathematics and Mathematical Sciences* 2003.21 (2003), pp. 1305–1321.

[40]    D. P. Kingma and J. Ba. "Adam: A Method for Stochastic Optimization". In: *arXiv preprint* 1412.6980 (2014).

[41]   E. Kobler et al. "Total Deep Variation for Linear Inverse Problems". In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition.* 2020, pp. 7549–7558.

[42]   T. Kohler, C. Bontus, and P. Koken. "The radon-split method for helical cone-beam CT and its application to nongated reconstruction". In: *IEEE Transactions on Medical Imaging* 25.7 (2006), pp. 882–897.

[43]   T. E. Komolafe et al. "Smoothed L0-Constraint Dictionary Learning for Low-Dose X-Ray CT Reconstruction". In: *IEEE Access* 8 (2020), pp. 116961–116973.

[44]   M. S. Lewicki and B. A. Olshausen. "Probabilistic framework for the adaptation and comparison of image codes". In: *JOSA A* 16.7 (1999), pp. 1587–1601.

[45]   M. S. Lewicki and T. J. Sejnowski. "Learning overcomplete representations". In: *Neural computation* 12.2 (2000), pp. 337–365.

[46]   Y. Liu and T. Kozubowski. "A folded Laplace distribution". In: *Journal of Statistical Distributions and Applications* 2 (2015), pp. 1–17.

[47]   S. Lunz, O. Öktem, and C.-B. Schönlieb. "Adversarial regularizers in inverse problems". In: *Advances in Neural Information Processing Systems.* 2018, pp. 8507–8516.

[48]   S. Mahdizadehaghdam et al. "Deep dictionary learning: A parametric network approach". In: *IEEE Transactions on Image Processing* 28.10 (2019), pp. 4790–4802.

[49]   C. H. McCollough et al. "Low-dose CT for the detection and classification of metastatic liver lesions: Results of the 2016 Low Dose CT Grand Challenge". In: *Medical Physics* 44.10 (2017), e339–e352.

[50]   C. H. McCollough et al. "Strategies for Reducing Radiation Dose in CT". In: *Radiology Clinics* 47.1 (2009), pp. 27–40.

[51]   A. Mileto et al. "State of the Art in Abdominal CT: The Limits of Iterative Reconstruction Algorithms". In: *Radiology* 293.3 (2019), pp. 491–503.

[52]   F. Natterer. "A Sobolev space analysis of picture reconstruction". In: *SIAM Journal on Applied Mathematics* 39.3 (1980), pp. 402–411.

[53]   F. Natterer and F. Wübbeling. *Mathematical Methods in Image Reconstruction.* Vol. 5. Monographs on Mathematical Modeling and Computation. Society for Industrial and Applied Mathematics, 2001. DOI: 10.1137/1.9780898718324.

[54]   Y. E. Nesterov. "A method of solving a convex programming problem with convergence rate $O(1/k^2)$". In: *Soviet Mathematics Doklady* 27.2 (1983), pp. 372–376.

[55]   F. Noo, J. Pack, and D. Heuscher. "Exact helical reconstruction using native cone-beam geometries". In: *Physics in Medicine & Biology* 48.23 (2003), pp. 3787–3818.

[56]   D. Obmann, J. Schwab, and M. Haltmeier. "Deep synthesis regularization of inverse problems". In: *arXiv preprint* 2002.00155 (2020).

[57] B. A. Olshausen and D. J. Field. "Sparse coding with an overcomplete basis set: A strategy employed by V1?" In: *Vision research* 37.23 (1997), pp. 3311–3325.

[58] X. Pan, E. Y. Sidky, and M. Vannier. "Why do commercial CT scanners still employ traditional, filtered back-projection for image reconstruction?" In: *Inverse Problems* 25.12 (2009), p. 123009.

[59] V. Papyan, Y. Romano, and M. Elad. "Convolutional neural networks analyzed via convolutional sparse coding". In: *The Journal of Machine Learning Research* 18.1 (2017), pp. 2887–2938.

[60] S. Ravishankar, J. C. Ye, and J. A. Fessler. "Image reconstruction: From sparsity to data-adaptive methods and machine learning". In: *Proceedings of the IEEE* 108.1 (2019), pp. 86–109.

[61] R. Rubinstein, A. M. Bruckstein, and M. Elad. "Dictionaries for Sparse Representation Modeling". In: *Proceedings of the IEEE* 98.6 (2010), pp. 1045–1057.

[62] L. I. Rudin, S. Osher, and E. Fatemi. "Nonlinear total variation based noise removal algorithms". In: *Physica D: nonlinear phenomena* 60.1-4 (1992), pp. 259–268.

[63] A. Sarma et al. "Radiation and Chest CT Scan Examinations: What Do We Know?" In: *Chest* 142.3 (2012), pp. 750–760.

[64] L. A. Shepp and B. F. Logan. "The Fourier reconstruction of a head section". In: *IEEE Transactions in Nuclear Science* 21.3 (1974), pp. 21–43.

[65] D. Simon and M. Elad. "Rethinking the CSC model for natural images". In: *Advances in Neural Information Processing Systems.* 2019, pp. 2274–2284.

[66] J. Sulam et al. "Multilayer convolutional sparse modeling: Pursuit and dictionary learning". In: *IEEE Transactions on Signal Processing* 66.15 (2018), pp. 4090–4104.

[67] J. Sulam et al. "On multi-layer basis pursuit, efficient algorithms and convolutional neural networks". In: *IEEE transactions on pattern analysis and machine intelligence* (2019).

[68] D. Tack and P. A. Gevenois, eds. *Radiation Dose from Adult and Pediatric Multidetector Computed Tomography-Springer.* Springer Verlag, 2007.

[69] A. M. Tillmann. "On the Computational Intractability of Exact and Approximate Dictionary Learning". In: *IEEE Signal Processing Letters* 22.1 (2015), pp. 45–49.

[70] B. Tolooshams, S. Dey, and D. Ba. "Deep Residual Autoencoders for Expectation Maximization-Inspired Dictionary Learning". In: *IEEE Transactions on Neural Networks and Learning Systems* (2020).

[71] J. A. Tropp and A. C. Gilbert. "Signal recovery from random measurements via orthogonal matching pursuit". In: *IEEE Transactions on information theory* 53.12 (2007), pp. 4655–4666.

[72] D. Ulyanov, A. Vedaldi, and V. Lempitsky. "Deep image prior". In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition.* 2018, pp. 9446–9454.

[73]  M. J. Willemink and P. B. Noël. "The evolution of image reconstruction for CT — from filtered back projection to artificial intelligence". In: *European Radiology* 29 (2019), pp. 2185–2195.

[74]  Q. Xu et al. "Low-dose X-ray CT reconstruction via dictionary learning". In: *IEEE transactions on medical imaging* 31.9 (2012), pp. 1682–1697.

[75]  H. Yu et al. "Studies on Palamodov's algorithm for cone-beam CT along a general curve". In: *Inverse Problems* 22.2 (2006), pp. 447–460.

[76]  C. Zhang et al. "Low-dose CT reconstruction via L1 dictionary learning regularization using iteratively reweighted least-squares". In: *Biomedical engineering online* 15.1 (2016), p. 66.

[77]  X. Zhao and J. Guo. "Low-dose CT Image Reconstruction via Total Variation and Dictionary Learning". In: *2019 IEEE 9th International Conference on Electronics Information and Emergency Communication (ICEIEC)*. IEEE. 2019, pp. 248–251.

[78]  J. Zhu et al. "Numerical studies on Feldkamp-type and Katsevich-type algorithms for cone-beam scanning along nonstandard spirals". In: *Developments in X-Ray Tomography IV, (26 October 2004)*. Ed. by U. Bonse. Vol. 5535. Proceedings of SPIE. 2004, 8 p.

## A Comparison of patch-based and a convolutional synthesis operators

We compare the performance of patch-based synthesis operator $\mathcal{S}_{\hat{\mathcal{D}}}^p$ and the convolutional operator $\mathcal{S}_{\hat{\mathcal{D}}}$ in the case of tomographic reconstruction problem eq. (36). This problem is reduced to eq. (37) in the case of convolutional synthesis operator. A similar formulation can be used when the linear synthesis operator $\mathcal{S}_{\hat{\mathcal{D}}}^p$ acts on patches:

$$(\widehat{x}, \widehat{z}) \in \underset{(x,z) \in X \times Z}{\arg\max} \, \mathcal{L}\big(\mathcal{A}(x), y\big) + \frac{\lambda_1}{N_p} \sum_{i=1}^{N_p} \big\|x_i - \mathcal{S}_{\hat{\mathcal{D}}}^p(z_i)\big\|_2^2 + \frac{\lambda_2}{N_p} \sum_{i=1}^{N_p} \|z_i\|_1 \quad (49)$$

Here $\{x_i\}_{i=1}^{N_p}, x_i \in \mathbb{R}^{k \times k}$ is a sequence of all overlapping $k \times k$ image patches, $\{z_i\}_{i=1}^{N_p} \subset \mathbb{R}^m$ is the sequence of corresponding coefficients with $m$ denoting the number of dictionary atoms. Finally, $N_p$ is the number of patches. In addition, we re-scale by a factor $k^{-2}$ simply to make values of regularization parameters $\lambda_1$ and $\lambda_2$ numerically closer to the values used in the convolutional setting eq. (37).

On our validation set (21 image slices) we observe PSNR values 47.94 and 48.22 for $\mathcal{S}_{\hat{\mathcal{D}}}^p$ and $\mathcal{S}_{\hat{\mathcal{D}}}$ respectively. In the case of patch-based synthesis, the best performance is achieved for $\lambda_1 = 10, \lambda_2 = 0.0006$. We show example reconstructions obtained with different values of regularization parameters in the second row of fig. 3. The best parameter setting in the patch-based case is highlighted in bold. Even in this setting, the image (fig. 3d) looks more noisy than the image synthesised with convolutional operator (fig. 3b).

We also tried to adapt the sparsity level during dictionary learning, to learn a better dictionary for the patch-based setting. However, this did not lead to significant improvements (best achieved PSNR is 48.00).

## B Learned dictionary

Figure 4 shows learned dictionary atoms ordered by significance. We measure the significance of each atom by summing absolute values of the corresponding dictionary coefficients. The coefficients are calculated by solving the reconstruction problem in eq. (37) for all images in the validation set.

Note that the learned dictionary also contains very general atoms that contain isolated point-like image features (bottom row in fig. 4). However, those atoms are used less significantly during reconstruction.

## C Quality of images depending on regularization parameters

We show the dependency of dictionary learning based regularization eq. (37) on regularization parameters $\lambda_1$ and $\lambda_2$ in fig. 5. The optimal setting $\lambda_1 = 50$ and $\lambda_2 = 0.0016$ is highlighted in bold.

The first parameter $\lambda_1$ controls the distance of reconstructed image $x$ from the image synthesised by the dictionary $\mathcal{S}_{\hat{\mathcal{D}}}(z)$. We show images for smaller values of $\lambda_1$ in the first row of fig. 5. We find that a wide range of values leads to similar results. For instance, for $\lambda_1 = 50$ and $\lambda_1 = 10$ images are very similar. Only for much lower values, like $\lambda_1 = 1$ we can see that the noise pattern changes

(a) True image

(b) Convolutional synthesis operator, $\lambda_1$=50, $\lambda_2$=0.0016, PSNR=44.49

(c) Patch-based synthesis operator, $\lambda_1$=10, $\lambda_2$=0.0016, PSNR=43.55

(d) **Patch-based synthesis operator, $\lambda_1$=10, $\lambda_2$=0.0006, PSNR=44.14**

(e) Patch-based synthesis operator, $\lambda_1$=10, $\lambda_2$=0.0003, PSNR=43.55
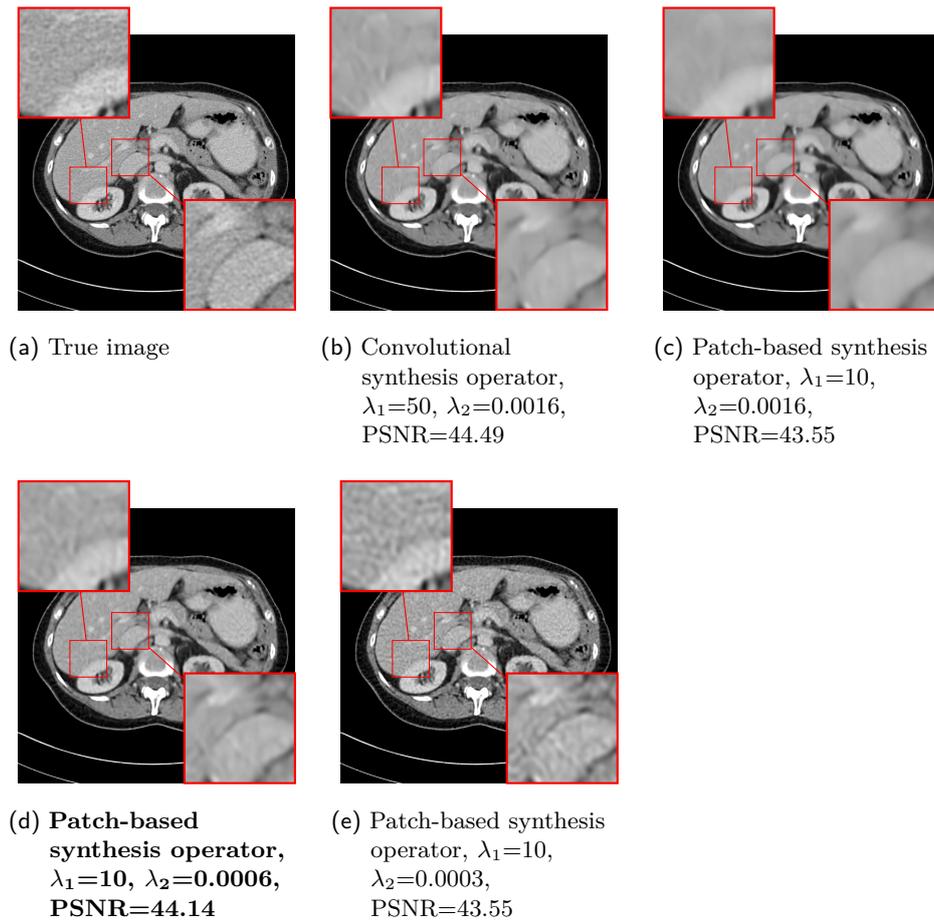
Fig. 3: Example slice of the abdomen (level 80 HU, window 370 HU) reconstructed with dictionary learning based regularization of image patches and different values of regularization parameters. Region of interests illustrate noise texture (top left) and delineation of low contrast structures (liver-kidney border in the top left region of interest and liver-pancreas transition as well as contrasted vessel in the liver in the bottom right region of interest.
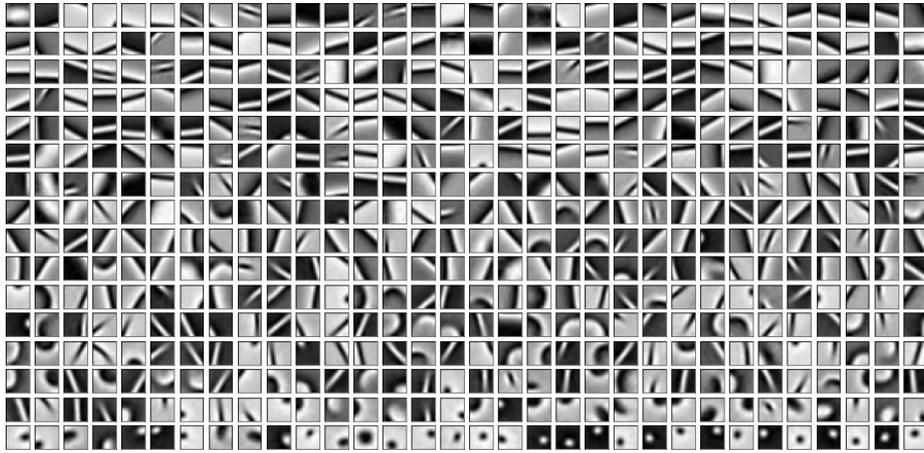
Fig. 4: Learned dictionary atoms ordered by significance (from top to bottom).

significantly. In particular, the low frequency noise shaped by the dictionary atoms becomes masked by high frequency noise. Unfortunately, quality of the image measured by PSNR significantly decreases.

The second parameter $\lambda_2$ controls sparsity of dictionary coefficients. We show images for different values of this parameter in the second row of fig. 5. We can see that for a higher value ($\lambda_2 = 0.0024$) image is over-smoothed and for a lower value ($\lambda_2 = 0.0012$) image is more noisy. Nevertheless, quality of the image measured by PSNR stays relatively high in both cases.

To summarize, the method seems robust to the choice of regularization parameter.

(a) True image

(b) $\lambda_1=10$, $\lambda_2=0.0016$, PSNR=44.24

(c) $\lambda_1=1$, $\lambda_2=0.0016$, PSNR=41.73

(d) $\lambda_1=50$, $\lambda_2=0.0024$, PSNR=44.37

(e) $\mathbf{\lambda_1=50}$, $\mathbf{\lambda_2=0.0016}$, **PSNR=44.49**

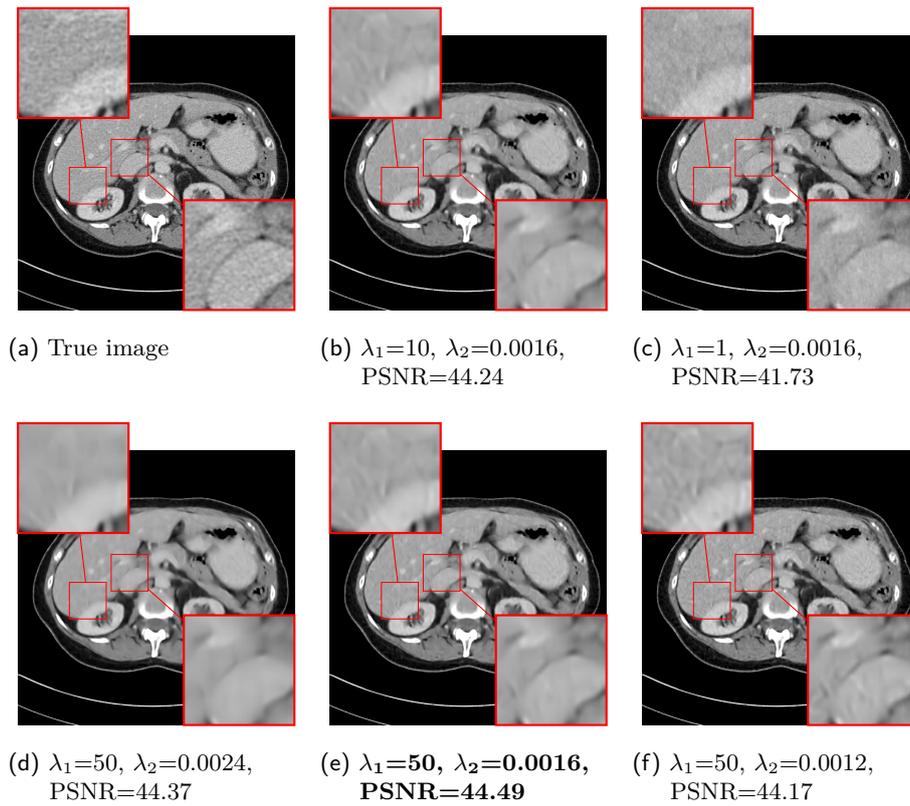(f) $\lambda_1=50$, $\lambda_2=0.0012$, PSNR=44.17

Fig. 5: Example slice of the abdomen (level 80 HU, window 370 HU) reconstructed with dictionary learning based regularization and different values of regularization parameters. Region of interests illustrate noise texture (top left) and delineation of low contrast structures (liver-kidney border in the top left region of interest and liver-pancreas transition as well as contrasted vessel in the liver in the bottom right region of interest.