

Deep Variational Networks with Exponential Weighting for Learning Computed Tomography

Valery Vishnevskiy¹, Richard Rau, and Orcun Goksel

Computer-assisted Applications in Medicine Group, ETH Zurich, Switzerland

¹valeryv@vision.ee.ethz.ch

Abstract. Tomographic image reconstruction is relevant for many medical imaging modalities including X-ray, ultrasound (US) computed tomography (CT) and photoacoustics, for which the access to full angular range tomographic projections might be not available in clinical practice due to physical or time constraints. Reconstruction from incomplete data in low signal-to-noise ratio regime is a challenging and ill-posed inverse problem that usually leads to unsatisfactory image quality. While informative image priors may be learned using generic deep neural network architectures, the artefacts caused by an ill-conditioned design matrix often have global spatial support and cannot be efficiently filtered out by means of convolutions. In this paper we propose to learn an inverse mapping in an end-to-end fashion via unrolling optimization iterations of a prototypical reconstruction algorithm. We herein introduce a network architecture that performs filtering jointly in both sinogram and spatial domains. To efficiently train such deep network we propose a novel regularization approach based on deep exponential weighting. Experiments on US and X-ray CT data show that our proposed method is qualitatively and quantitatively superior to conventional non-linear reconstruction methods as well as state-of-the-art deep networks for image reconstruction. Fast inference time of the proposed algorithm allows for sophisticated reconstructions in real-time critical settings, demonstrated with US SoS imaging of an *ex vivo* bovine phantom.

1 Introduction

Tomographic image reconstruction with sparse or limited angular (LA) data arises in a number of applications including image guided interventions [13], photoacoustics [7], and US speed-of-sound (SoS) imaging [11,2]. Such underdetermined problems usually require suitable problem-specific regularization for meaningful reconstructions, e.g. free from streaking artefacts. Setting regularization parameters manually can be cumbersome and often generalizes poorly. Using learning-based methods as in [18] can greatly improve reconstruction accuracy and account for non-Gaussian noise models. Unfortunately the method in [18] is based on patch-based clustering leading to very slow reconstruction. Straightforward application of computationally efficient convolutional network directly to measurements might be tempting, but is unjustified, because sinogram values have global spatial dependence on image intensities. In practice, such

generic networks are not likely to generalize well for LA-CT problems [9]. Many deep-learning inspired methods employ artificial neural networks to learn filtering or weighting [17,12] of the input sinograms prior to the backprojection step, after which the result might be postprocessed by another network [4]. Unfortunately such sinogram preprocessing requires problem-specific weighting schemes, which would constrain applicable acquisition geometries. Variational Networks (VN) employ adjoint of projection operator to learn convolutional filters in *spatial* domain. For compressed sensing in MRI, a landmark VN architecture was introduced by Hammernik et al. [3], which in practice relies on unitarity of Fourier transform. This was addressed in [15] with more sophisticated unrolled iterations that improved USCT reconstructions and allowed *detection* of coarse blob-looking inclusions.

In order to allow for accurate image reconstruction with ill-conditioned spatial encoding operators, in this paper we extend VN architecture by introducing sinogram filters that are learned as preconditioners, inspired by filtered backprojection. We also propose an efficient network regularization scheme that allows stable training inspired by Landweber iterations [6] and deep supervision [8].

2 Methods

Tomographic reconstruction problem involves estimating image intensities x from set of measurements (sinogram) b_i (e.g., time-of-flight or ray attenuation) that are modelled, e.g., as line integrals $b_i = \int_{\text{ray}_i} x(\mathbf{r})ds$. Given a set of measurements $\mathbf{b} \in \mathbb{R}^M$, algebraic reconstruction methods solve for discretized spatial encoding equations in a maximum-a-posteriori (MAP) sense, i.e.:

$$\hat{\mathbf{x}}(\mathbf{b}; \lambda, p) = \underset{\mathbf{x}}{\operatorname{argmin}} \|\mathbf{L}\mathbf{x} - \mathbf{b}\|_p^p + \lambda\mathcal{R}(\mathbf{x}), \quad (1)$$

where $\mathbf{x} \in \mathbb{R}^{n_1 n_2}$ is $n_1 \times n_2$ image and $\mathbf{L} \in \mathbb{R}^{M \times n_1 n_2}$ is a sparse, ray-discretization matrix that depends on the acquisition geometry. The norm p determines data inconsistency penalty, e.g. $p=2$ assumes Gaussian acquisition noise while $p=1$ assumes Laplace noise, the latter of which is often considered to be more robust. Non-negative weight λ controls the influence of regularizer. Similarly to [5], we herein consider total variation $\mathcal{R}_{\text{TV}}(\mathbf{x}) = \|\nabla\mathbf{x}\|_1$ and total generalized variation $\mathcal{R}_{\text{TGV}}(\mathbf{x}) = \min_{\mathbf{u}} \|\nabla\mathbf{x} - \mathbf{u}\|_1 + 2\|\mathcal{E}\mathbf{u}\|_1$ regularizers. Here ∇ denotes first-order forward finite derivative matrix and \mathcal{E} is the symmetrized vector field derivative operator. Both TV and TGV regularizers yield convex optimization problems for $p = \{1, 2\}$, which we hereafter refer as $L_p\text{TV}$ and $L_p\text{TGV}$.

A Variational Network can be seen as a sequence of K unrolled iterations of a numerical optimization scheme, inspired by a prototypical objective as in (1). For additional learning capacity, these iterations are further relaxed, e.g. by adding variable filters and activations that can be tuned during training. As illustrated in Fig. 1, we initialize VN inference by $\mathbf{L}^T\mathbf{b}$. Then, at each network layer (iteration) k , a *gradient* term $\mathbf{g}^{(k)}$ is accumulated via running average governed by

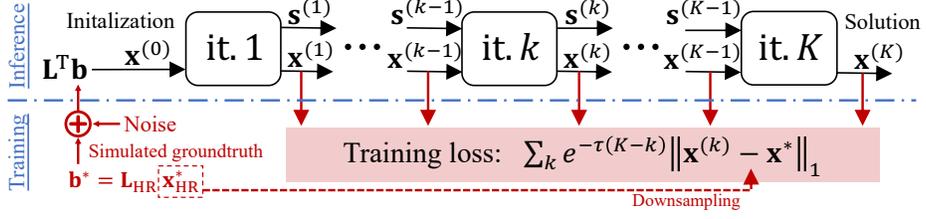


Fig. 1. Structure of variational network and its training strategy.

momentum coefficient $\alpha^{(k)}$ and used to update current image estimate $\mathbf{x}^{(k)}$ as follows:

$$\begin{aligned} \mathbf{g}^{(k)} &\leftarrow \left(\mathbf{P}^{(k)} \mathbf{L} \mathbf{Q}^{(k)} \right)^{\top} \mathbf{W}_d^{(k)} \varphi_d^{(k)} \left\{ \mathbf{W}_d^{(k)} \left(\mathbf{P}^{(k)} \mathbf{L} \mathbf{Q}^{(k)} \mathbf{x}^{(k-1)} - \mathbf{b} \right) \right\} + \\ &\quad \left(\mathbf{D}^{(k)} \right)^{\top} \mathbf{W}_r^{(k)} \varphi_r^{(k)} \left\{ \mathbf{W}_r^{(k)} \mathbf{D}^{(k)} \mathbf{x}^{(k-1)} \right\}, \\ \mathbf{s}^{(k)} &\leftarrow \alpha^{(k)} \mathbf{s}^{(k-1)} + \mathbf{g}^{(k)}, \quad \mathbf{x}^{(k)} \leftarrow \mathbf{x}^{(k-1)} - \mathbf{s}^{(k)}. \end{aligned} \quad (2)$$

We define the gradient term $\mathbf{g}^{(k)}$ with the following tunable operations that all allow backpropagation of gradients: (i) multiplication with diagonal *pre-conditioner* $\mathbf{W}_d^{(k)}$ and spatial regularization *weighting* $\mathbf{W}_r^{(k)}$; (ii) convolution with left ($\mathbf{P}^{(k)}$) and right ($\mathbf{Q}^{(k)}$) preconditioners, and several (herein, $n_f=50$) regularization filters $\mathbf{D}^{(k)}$; (iii) nonlinear data ($\varphi_d^{(k)}$) and regularization ($\varphi_r^{(k)}$) activation functions that are parametrized via linear interpolation on a regular grid, herein, of size $n_g=35$; i.e. $\varphi\{t\} = (1 - t + \lfloor t \rfloor) \phi_{\lfloor t \rfloor} + (t - \lfloor t \rfloor) \phi_{\lfloor t \rfloor + 1}$. To avoid bilinear ambiguities, every $n_k \times n_k$ (herein, $n_k=7$) convolution \mathbf{D} , \mathbf{Q} , \mathbf{P} with kernel \mathbf{d} is reparametrized to be zero-centered unit-norm, i.e. $\mathbf{d} = n_k(\mathbf{d}' - \text{mean}(\mathbf{d}')) / \text{std}(\mathbf{d}')$, while diagonal terms are also ensured to be nonnegative and bounded via sigmoid: $\mathbf{W} = \text{diag}(\sigma(w_i))$. Stochastic minimization of the *exponentially weighted* ℓ_1 reconstruction loss is then conducted to tune parameter set $\Theta = \{\mathbf{P}^{(k)}, \mathbf{Q}^{(k)}, \mathbf{D}^{(k)}, \mathbf{W}_r^{(k)}, \mathbf{W}_f^{(k)}, \varphi_r^{(k)}, \varphi_d^{(k)}, \alpha^{(k)}\}$:

$$\min_{\Theta} \mathbb{E}_{\{\mathbf{b}, \mathbf{x}^*\} \in \mathcal{T}} \sum_{k=1}^K e^{-\tau(K-k)} \|\mathbf{x}^{(k)}(\mathbf{b}; \Theta) - \mathbf{x}^*\|_1, \quad (3)$$

on the training set \mathcal{T} . Here $\tau \geq 0$ controls the regularization of the network: at $\tau=0$, the reconstruction on all layers is weighted equally, therefore all network parameters have low variance of gradients, which allows for stable training. For $\tau \rightarrow +\infty$, only the last network output $\mathbf{x}^{(K)}$ is used for training, which allows the network to be tuned accurately for the final objective. We accordingly increase τ during the training procedure to gradually relax constraints on the network. Intuitively, such regularization encourages VN to provide reconstruction as early as possible, which is inspired by *early stopping* — a common image reconstruction strategy that allows to avoid degenerate solutions and can be shown to be equivalent to Tikhonov regularization in certain cases [6].

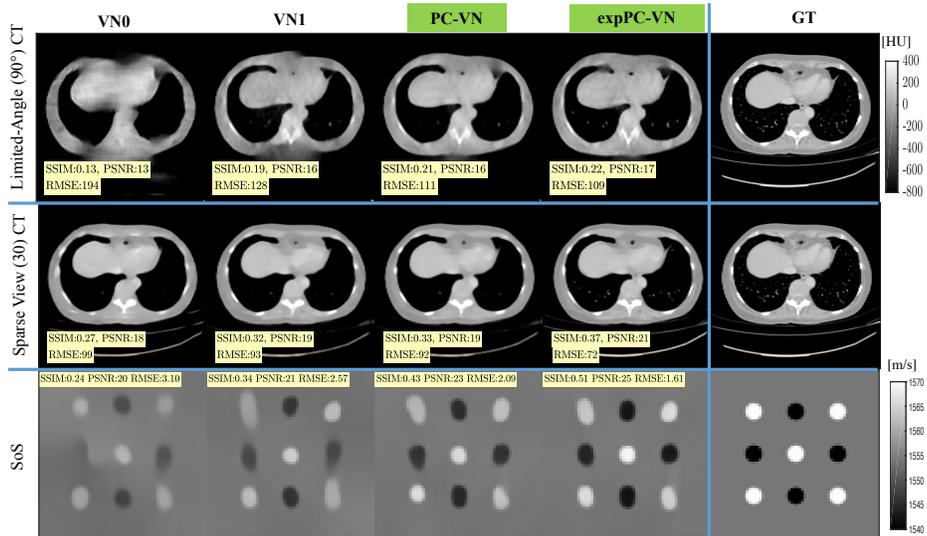


Fig. 2. Performance comparison of VN architectures on X-ray CT and SoS simulations. Proposed architectures are highlighted in green.

Training. We employ a $K=10$ layer VN and perform $5 \cdot 10^4$ iterations of Adam algorithm (learning rate 10^{-3} , $\beta_1=0.85$, $\beta_2=0.98$, batch size of 10) for training, during which we continually adjust $\tau = j \cdot 10^{-3}$ with j being the iteration number. To avoid overfitting to the discretization scheme employed in the simulation, we use higher resolution images \mathbf{x}_{HR} to compute line integrals defined by \mathbf{L}_{HR} while training for these to be reconstructed in the desired resolution defined by the discretization of \mathbf{L} . The network was implemented using Tensorflow framework, where multiplication with the design matrix \mathbf{L} was carried out as a generic sparse matrix-vector multiplication. For comparison with iterative reconstruction, we employ ADMM algorithm [1] to solve (1) by approximating the regularized inversion of \mathbf{L} with 5 iterations of LSQR solver [10]. For each experimental scenario, an optimal value of λ was tuned on a single test image.

3 Experiments and Results

X-ray CT dataset. We used 3DIRCADb dataset [14] that consists of 3080 axial CT scans from 22 patients with inplane resolution varying from 0.56^2 to 0.961^2 mm^2 and slice thickness of 1.6 to 4 mm. The forward simulation was conducted on original 512×512 images, which then were downsampled to 256×256 grid, yielding the ground truth (GT). We simulated parallel beam acquisition geometry for limiter-angle (LA) and sparse-view (SV) scenarios. In the LA scenario we simulated angular ranges of 120° , 90° and 60° , where projections were acquired in 1° increments. For the SV scenario we simulated 180° range with 60, 30, and 15 uniformly-acquired projections. To simulate realistic acquisition noise, we

Table 1. Mean reconstruction RMSE computed on corresponding training sets with standard deviations indicated in parentheses. Proposed methods are highlighted.

Dataset	L2TV	L1TV	L2TGV	VN0	VN1	PC-VN	expPC-VN
	RMSE	RMSE	RMSE	RMSE	RMSE	RMSE	RMSE
LA-CT 90°	216 (9)	238 (16)	205 (6)	210 (9)	128 (15)	119 (19)	103 (8)
SV-CT 30	93 (5)	91 (4)	98 (6)	97 (4)	87 (8)	88 (9)	65 (5)
SoS USCT	1.48 (0.44)	1.62 (0.48)	1.89 (0.35)	1.8 (0.33)	1.35 (0.43)	1.19 (0.41)	1.0 (0.35)

follow [18] and employ Poisson+Gaussian model, i.e. $b = \log |\text{Poisson}(I_0 \exp(-b^*)) + \text{Gauss}(0, \sigma_E)|$, with $I_0=2 \cdot 10^4$ and $\sigma_E=8 \cdot \text{Unif}(0, 1)$ to allow variable SNR. We used 20 patient scans for training and two for testing.

US Speed-of-Sound Tomography. We follow [15] to simulate reflector-based USCT reconstruction with a 128-element transducer and a square imaging field-of-view. Synthetic inclusion masks were generated at 256×256 resolution as levelsets of random, spatially-smooth functions. The inclusion SoS values were then randomly sampled from $[1350, 1650]$ m/s. Acquisition noise was modelled as Gaussian with $\sigma_N=2 \cdot 10^{-8}$ and the reconstruction was conducted on a coarser 64×64 grid to avoid overfitting to the discretization. The training set contains 15000 random synthetic inclusions, while the test set includes 13 geometric primitives consisting of oval and polygonal inclusions.

Evaluation. For quantitative comparison of reconstruction and ground truth, we calculated structural similarity index measurement (SSIM) [16], Root Mean Square Error (RMSE), and peak signal-to-noise ratio (PSNR) as follows:

$$\text{RMSE}(\mathbf{x}, \mathbf{y}) = \sqrt{\frac{\|\mathbf{x} - \mathbf{y}\|_2^2}{N}}, \quad \text{PSNR}(\mathbf{x}, \mathbf{y}) = 10 \log_{10} \frac{R^2 N}{\|\mathbf{x} - \mathbf{y}\|_2^2}, \quad (4)$$

where R is the dynamic range of the ground truth image. The corresponding values are reported in Hounsfield units (HU) and m/s, for X-ray CT and USCT SoS experiments accordingly. We denote the architectures employed in [3] and [15] as VN0 and VN1, respectively. As seen in Fig. 2 the proposed preconditioned network (PC-VN) with sinogram convolutions improves reconstruction accuracy and quality for all X-ray CT acquisition scenarios and USCT as suggested by RMSE and SSIM values. As reported in Tab. 1, training the proposed network using exponentially weighted loss (expPC-VN) defined in Eq. (3) further improves reconstruction quality and reduces the variance of error, which can be explained by the introduced regularization effect. Fig. 3 shows that expPC-VN outperforms iterative methods both in terms of accuracy and image quality. Namely, compared to nonlinear reconstruction methods, we observe improvements of RMSE by 49% in the CT-LA-60° scenario, and increase of SSIM by 38% in the CT-SV-15 experiment. Quantitative results from Tab. 1 and Fig. 3 show that proposed reconstruction method outperforms all considered iterative and deep learning -based approaches.

In Fig. 4(a), we present a USCT SoS reconstruction from *ex vivo* bovine skeletal muscle tissue embedded in a gelatin phantom. Compared to the conven-

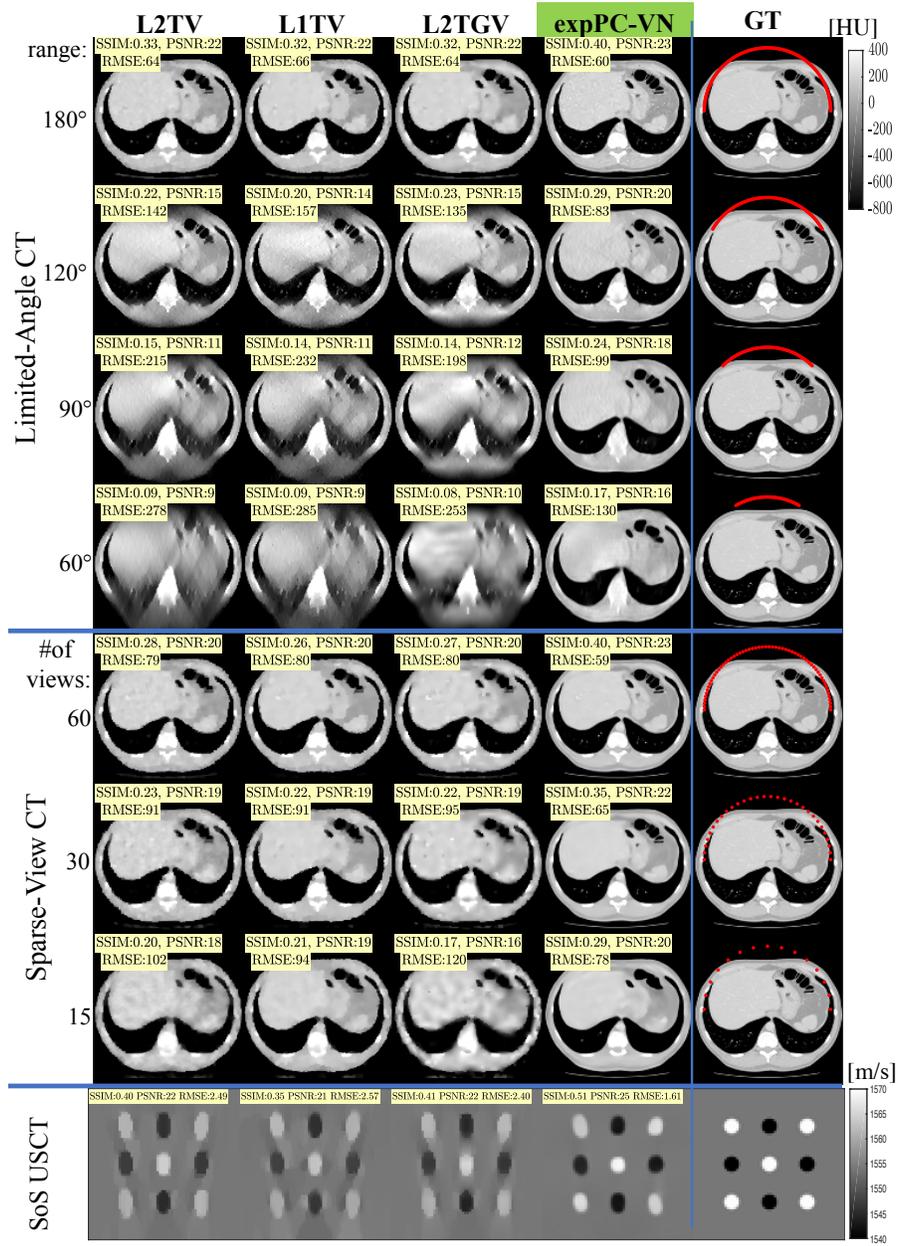


Fig. 3. Reconstruction results for X-ray CT and USCT acquisitions. For sparse view and limited angle experiments, acquired angular positions of projections are depicted in the GT column.

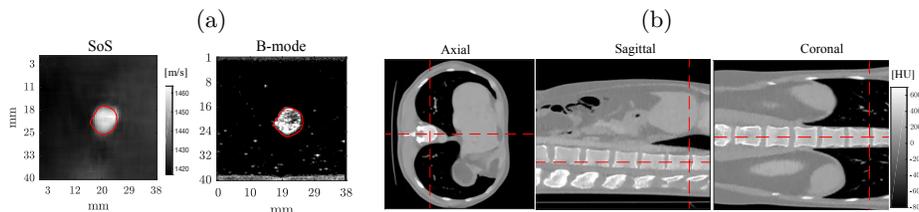


Fig. 4. (a) expPC-VN SoS reconstruction of *ex vivo* phantom and corresponding B-mode image. Red contour shows inclusion segmentation from the B-mode image. (b) Axial, sagittal and coronal view of slice-wise reconstruction of the proposed expPC-VN reconstruction network. Corresponding slice positions are indicated with red dashed line.

tional B-Mode image, we could accurately identify inclusion location and provide quantitative estimates of local tissue SoS. Reconstruction of a single image with VN takes 0.03 s on NVIDIA Titan Xp GPU and 1-4 min with iterative methods on a 6-core 3.7 GHz Intel CPU. In order to demonstrate potential 3D imaging applications of our method, we also conducted X-ray CT reconstruction of SV-60 acquisition scenario with test images rotated by 90° in the axial plane, and show a cross-sectional views from the reconstructed 3D volume in Fig. 4 (b). We observe high spatial coherence and contrast in coronal and sagittal planes which asserts high generalization ability of the proposed expPC-VN method.

4 Conclusions

In this paper we have presented a network architecture for preconditioned reconstruction and a regularization scheme for its efficient training via exponential weighting. The proposed network has been shown to outperform conventional algebraic and learning-based reconstruction methods in terms of accuracy and image quality for various challenging X-ray CT and SoS USCT scenarios. Such effectiveness and versatility of our approach may suggest its potential for solving other intriguing optimization and inverse problems also outside of the image reconstruction field.

References

1. Boyd, S., Parikh, N., Chu, E., Peleato, B., Eckstein, J.: Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends in ML* **3**(1), 1–122 (2011)
2. Cheng, A., Kim, Y., Anas, E.M., Rahmim, A., Boctor, E.M., Seifabadi, R., Wood, B.J.: Deep learning image reconstruction method for limited-angle ultrasound tomography in prostate cancer. In: *Procs SPIE Medical Imaging*. p. 1095516 (2019)
3. Hammernik, K., Klatzer, T., Kobler, E., Recht, M.P., Sodickson, D.K., Pock, T., Knoll, F.: Learning a variational network for reconstruction of accelerated MRI data. *MRM* **79**(6), 3055–3071 (2018)

4. Hammernik, K., Würfl, T., Pock, T., Maier, A.: A deep learning architecture for limited-angle computed tomography reconstruction. In: *Bildverarbeitung für die Medizin 2017*, pp. 92–97 (2017)
5. Knoll, F., Bredies, K., Pock, T., Stollberger, R.: Second order total generalized variation (TGV) for MRI. *MRM* **65**(2), 480–491 (2011)
6. Landweber, L.: An iteration formula for Fredholm integral equations of the first kind. *American journal of mathematics* **73**(3), 615–624 (1951)
7. Lin, H., Azuma, T., Unlu, M.B., Takagi, S.: Evaluation of adjoint methods in photoacoustic tomography with under-sampled sensors. In: *MICCAI*. pp. 73–81 (2018)
8. Liu, Y., Lew, M.S.: Learning relaxed deep supervision for better edge detection. In: *CVPR*. pp. 231–240 (2016)
9. Maier, A., Syben, C., Lasser, T., Riess, C.: A gentle introduction to deep learning in medical image processing. *Zeitschrift für Medizinische Physik* (2019)
10. Paige, C.C., Saunders, M.A.: LSQR: An algorithm for sparse linear equations and sparse least squares. *ACM TOMS* **8**(1), 43–71 (1982)
11. Sanabria, S.J., Goksel, O.: Hand-held sound-speed imaging based on ultrasound reflector delineation. In: *MICCAI*. pp. 568–576 (2016)
12. Schwab, J., Antholzer, S., Haltmeier, M.: Learned backprojection for sparse and limited view photoacoustic tomography. In: *Procs SPIE Photons Plus Ultrasound: Imaging and Sensing*. p. 1087837 (2019)
13. Siewerdsen, J., Daly, M., Bachar, G., Moseley, D., Bootsma, G., Brock, K., Ansell, S., Wilson, G., Chhabra, S., Jaffray, D., et al.: Multimode C-arm fluoroscopy, tomosynthesis, and cone-beam CT for image-guided interventions: from proof of principle to patient protocols. In: *Procs SPIE Medical Imaging*. p. 65101A (2007)
14. Soler, L., Hostettler, A., Agnus, V., Charnoz, A., Fasquel, J., Moreau, J., Osswald, A., Bouhadjar, M., Marescaux, J.: 3D image reconstruction for comparison of algorithm database: A patient specific anatomical and medical image database. *IRCAD, Strasbourg, France, Tech. Rep* (2010)
15. Vishnevskiy, V., Sanabria, S.J., Goksel, O.: Image reconstruction via variational network for real-time hand-held sound-speed imaging. *MLMIR* pp. 120–128 (2018)
16. Wang, Z., Bovik, A.C., Sheikh, H.R., Simoncelli, E.P., et al.: Image quality assessment: from error visibility to structural similarity. *IEEE Trans Imag Proc* **13**(4), 600–612 (2004)
17. Würfl, T., Hoffmann, M., Christlein, V., Breininger, K., Huang, Y., Unberath, M., Maier, A.K.: Deep learning computed tomography: learning projection-domain weights from image domain in limited angle problems. *IEEE Trans Med Imag* **37**(6), 1454–1463 (2018)
18. Zheng, X., Ravishankar, S., Long, Y., Fessler, J.A.: PWLS-ULTRA: An efficient clustering and learning-based approach for low-dose 3D CT image reconstruction. *IEEE Trans Med Imag* **37**(6), 1498–1510 (2018)