

Convergent Deep Convex and Non-Convex Regularization

Srijith Nair and Philip Schniter

Supported by NSF grant CCF-1955587



THE OHIO STATE UNIVERSITY

John D. and Alice Nelson Kraus Memorial
ECE Graduate Student Poster Competition - 2023

Variational signal/image recovery

Goal: Recover signal/image \mathbf{x}_0 from noisy measurements $\mathbf{y} = \mathcal{F}(\mathbf{x}_0, \mathbf{w})$

- Applications: denoising, deblurring, superresolution, inpainting, computed tomography, magnetic resonance imaging (MRI), phase-retrieval, de-quantization, photon-limited image recovery, etc.

Variational formulation:

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} J(\mathbf{x}) \text{ with } J(\mathbf{x}) \triangleq d(\mathbf{x}; \mathbf{y}) + r(\mathbf{x}; \boldsymbol{\theta})$$

- $d(\cdot; \mathbf{y})$ is a **data-fidelity term**, usually chosen as negative log-likelihood
- $r(\cdot; \boldsymbol{\theta})$ is a **regularizer** that incorporates prior information
- When r is differentiable, a prototypical algorithm is FBS/PGD:

$$\mathbf{x}^{(k+1)} = \text{prox}_{\tau^{(k)}d}(\mathbf{x}^{(k)} - \tau^{(k)}\nabla r(\mathbf{x}^{(k)}; \boldsymbol{\theta})),$$

where **explicit r** allows the use of back-tracking line-search to adapt $\tau^{(k)}$

- Key question: How do we choose the regularizer r ?

Regularizer design:

- Traditional hand-crafted regularizers use total-variation or wavelet sparsity
- We focus on **data-driven regularization** leveraging training data $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$

Data-driven regularization

Overall goals:

- Variational optimization should be **tractable** (e.g., J convex)
- Resulting $\hat{\mathbf{x}}$ should be close to ground truth \mathbf{x}_0

Regularizer architecture:

- Many have been proposed, based on, e.g., one-layer networks [1], denoisers [2, 3], autoencoders [4], and deep convolutional networks [5, 6, 7]
- Possible to make $r(\cdot; \boldsymbol{\theta})$ **structurally convex** [8, 9], but this requires specialized architectures that can limit expressivity

Gradient-step (GS) optimization [6, 10]:

- Train $I - \nabla r$ as a denoiser \Leftrightarrow train ∇r as a noise-estimator:

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \mathcal{L}_{\text{GS}}(\boldsymbol{\theta}) \quad \text{for} \quad \mathcal{L}_{\text{GS}}(\boldsymbol{\theta}) \triangleq \sum_{i=1}^n \mathbb{E}[L(\nabla r(\mathbf{x}_i + \mathbf{n}; \boldsymbol{\theta}), \mathbf{n})]$$

with $L(\hat{\mathbf{x}}, \mathbf{x})$ a loss like $L_2(\hat{\mathbf{x}}, \mathbf{x}) \triangleq \|\hat{\mathbf{x}} - \mathbf{x}\|_2^2$,

- Use automatic differentiation to compute ∇r
- Minimize $\tilde{J}(\mathbf{x}) = d(\mathbf{x}; \mathbf{y}) + \lambda r(\mathbf{x}; \boldsymbol{\theta})$ with λ tuned using validation data

Bi-level (BL) optimization [11]:

- Choose $\boldsymbol{\theta}$ to give a accurate $\hat{\mathbf{x}}_i(\boldsymbol{\theta}) \triangleq \arg \min_{\mathbf{x}} J(\mathbf{x}; \mathbf{y}_i, \boldsymbol{\theta})$, i.e.,

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \mathcal{L}_{\text{BL}}(\boldsymbol{\theta}) \quad \text{for} \quad \mathcal{L}_{\text{BL}}(\boldsymbol{\theta}) \triangleq \sum_{i=1}^n L(\hat{\mathbf{x}}_i(\boldsymbol{\theta}), \mathbf{x}_i)$$

- Can show that gradient equals

$$\nabla \mathcal{L}_{\text{BL}}(\boldsymbol{\theta}) = \sum_{i=1}^n \nabla_{\boldsymbol{\theta}} \nabla_{\mathbf{x}}^{\top} J(\hat{\mathbf{x}}_i; \mathbf{y}_i, \boldsymbol{\theta}) \boldsymbol{\gamma}_i$$

for the $\boldsymbol{\gamma}_i$ that solves the linear system

$$[\nabla_{\mathbf{x}} \nabla_{\mathbf{x}}^{\top} J(\hat{\mathbf{x}}_i; \mathbf{y}_i, \boldsymbol{\theta})] \boldsymbol{\gamma}_i = \mathbf{x}_i - \hat{\mathbf{x}}_i$$

- Use PGD to compute $\hat{\mathbf{x}}_i$, MINRES to compute $\hat{\boldsymbol{\gamma}}_i$, then gradient step in $\boldsymbol{\theta}$

Monotonicity

- A function $g : \mathbb{R}^p \rightarrow \mathbb{R}$ is convex if and only if its gradient is **monotone**, i.e.,

$$\frac{[\partial g(\mathbf{x}_1) - \partial g(\mathbf{x}_2)]^{\top} (\mathbf{x}_1 - \mathbf{x}_2)}{\|\mathbf{x}_1 - \mathbf{x}_2\|_2^2} \geq 0 \quad \forall \mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^p$$

where ∂g is the subdifferential

- Pesquet et al. [12] designed a data-driven denoiser that is the resolvent of a monotone operator, ensuring a convergent **PnP** algorithm [13]
- Below, we propose two ways to exploit monotonicity with an *explicit differentiable* regularizer r that has more in common with **RED** [2] than PnP

Proposed convex regularizer

- We propose to design r using an **adversarial monotonicity penalty on ∇r**

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \{ \mathcal{L}_{*}(\boldsymbol{\theta}) + \kappa_{\text{mon}} \mathcal{P}_{\nabla r}^{\text{mon}}(\boldsymbol{\theta}) + \kappa_{\text{lip}} \mathcal{P}_{\nabla r}^{\text{lip}}(\boldsymbol{\theta}) \}$$

$$\mathcal{P}_{\nabla r}^{\text{mon}}(\boldsymbol{\theta}) \triangleq \sum_{i=1}^n \max\{0, -\hat{\mu}_r(\mathbf{x}_i, \boldsymbol{\theta})\}$$

$$\mathcal{P}_{\nabla r}^{\text{lip}}(\boldsymbol{\theta}) \triangleq \sum_{i=1}^n \max\{0, \hat{\beta}_r(\mathbf{x}_i, \boldsymbol{\theta}) - \beta_r\}$$

$$\hat{\mu}_r(\mathbf{x}_i, \boldsymbol{\theta}) \triangleq \min_{\boldsymbol{\delta}^1, \boldsymbol{\delta}^2 \in \mathcal{B}(\epsilon)} \frac{[\nabla r(\mathbf{x}_i + \boldsymbol{\delta}^1; \boldsymbol{\theta}) - \nabla r(\mathbf{x}_i + \boldsymbol{\delta}^2; \boldsymbol{\theta})]^{\top} (\boldsymbol{\delta}^1 - \boldsymbol{\delta}^2)}{\|\boldsymbol{\delta}^1 - \boldsymbol{\delta}^2\|_2^2}$$

$$\hat{\beta}_r(\mathbf{x}_i, \boldsymbol{\theta}) \triangleq \max_{\boldsymbol{\delta}^3, \boldsymbol{\delta}^4 \in \mathcal{B}(\epsilon)} \frac{\|\nabla r(\mathbf{x}_i + \boldsymbol{\delta}^3; \boldsymbol{\theta}) - \nabla r(\mathbf{x}_i + \boldsymbol{\delta}^4; \boldsymbol{\theta})\|_2}{\|\boldsymbol{\delta}^3 - \boldsymbol{\delta}^4\|_2}$$

where $\mathcal{L}_{*} \in \{\mathcal{L}_{\text{BL}}, \mathcal{L}_{\text{GS}}\}$ and $\mathcal{B}(\epsilon)$ is an ϵ -ball with suitable ϵ

- The **Lipschitz penalty on ∇r** helps to stabilize training and improve performance

Proposed non-convex regularizer

- We propose to design r using an **adversarial monotonicity penalty on ∇J**

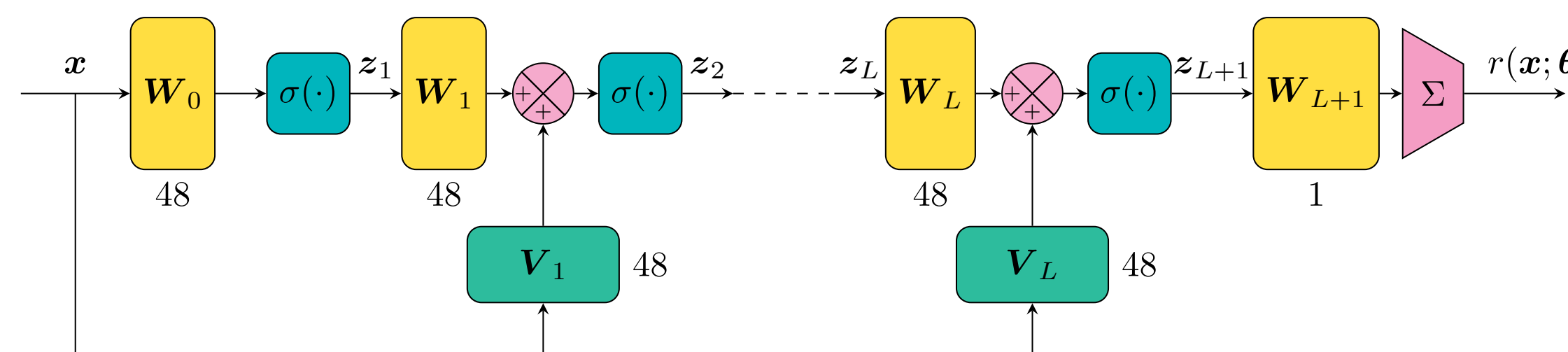
$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \{ \mathcal{L}_{\text{BL}}(\boldsymbol{\theta}) + \kappa_{\text{mon}} \mathcal{P}_{\nabla J}^{\text{mon}}(\boldsymbol{\theta}) + \kappa_{\text{lip}} \mathcal{P}_{\nabla r}^{\text{lip}}(\boldsymbol{\theta}) \}$$

$$\mathcal{P}_{\nabla J}^{\text{mon}}(\boldsymbol{\theta}) \triangleq \sum_{i=1}^n \max\{0, -\hat{\mu}_J(\mathbf{x}_i, \boldsymbol{\theta})\}$$

$$\hat{\mu}_J(\mathbf{x}_i, \boldsymbol{\theta}) \triangleq \min_{\boldsymbol{\delta}^1, \boldsymbol{\delta}^2 \in \mathcal{B}(\epsilon)} \frac{[\nabla J(\mathbf{x}_i + \boldsymbol{\delta}^1; \boldsymbol{\theta}) - \nabla J(\mathbf{x}_i + \boldsymbol{\delta}^2; \boldsymbol{\theta})]^{\top} (\boldsymbol{\delta}^1 - \boldsymbol{\delta}^2)}{\|\boldsymbol{\delta}^1 - \boldsymbol{\delta}^2\|_2^2}$$

- This encourages J to be convex without forcing r to be convex
- Reminiscent of the “convex non-convex” approach for handcrafted sparse regularizers [14]

Example regularizer architecture



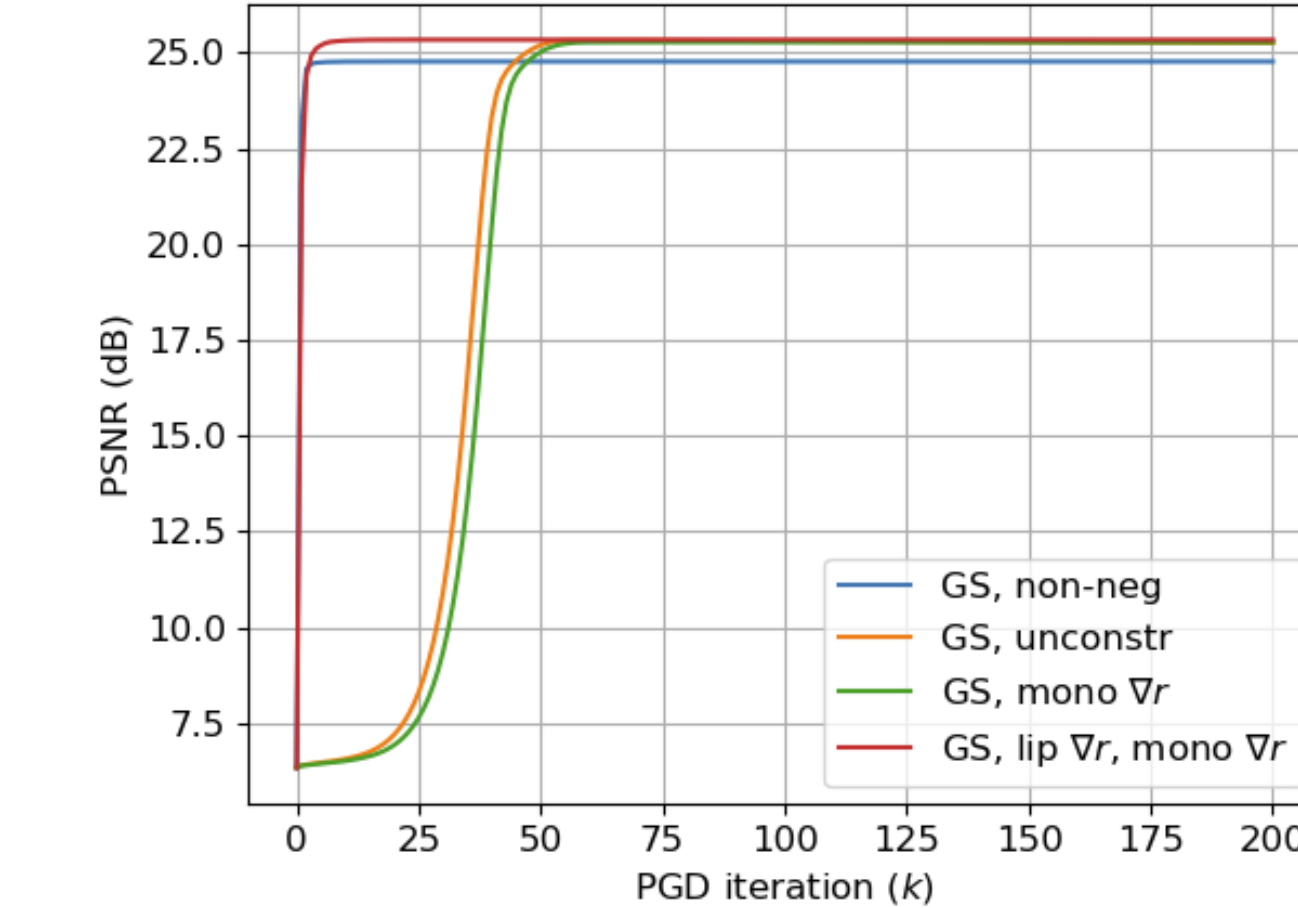
- Structurally convex under non-negativity constraint $[\mathbf{W}_l]_{jk} \geq 0$ [8]
 - Facilitates a direct comparison to our monotone approach
- Like Cohen et al. [6], we use it without constraints
- We use SiLU activations [15] (for differentiability) and $L = 7$ layers

Gaussian deblurring experiments (preliminary results)

- 25x25 Gaussian blur kernel with stdv 1.6
- Gaussian noise with stdv 10/255
- BSD400 dataset: 180x180 images, 350 training, 50 validation, 68 test

Gradient-step (GS) training

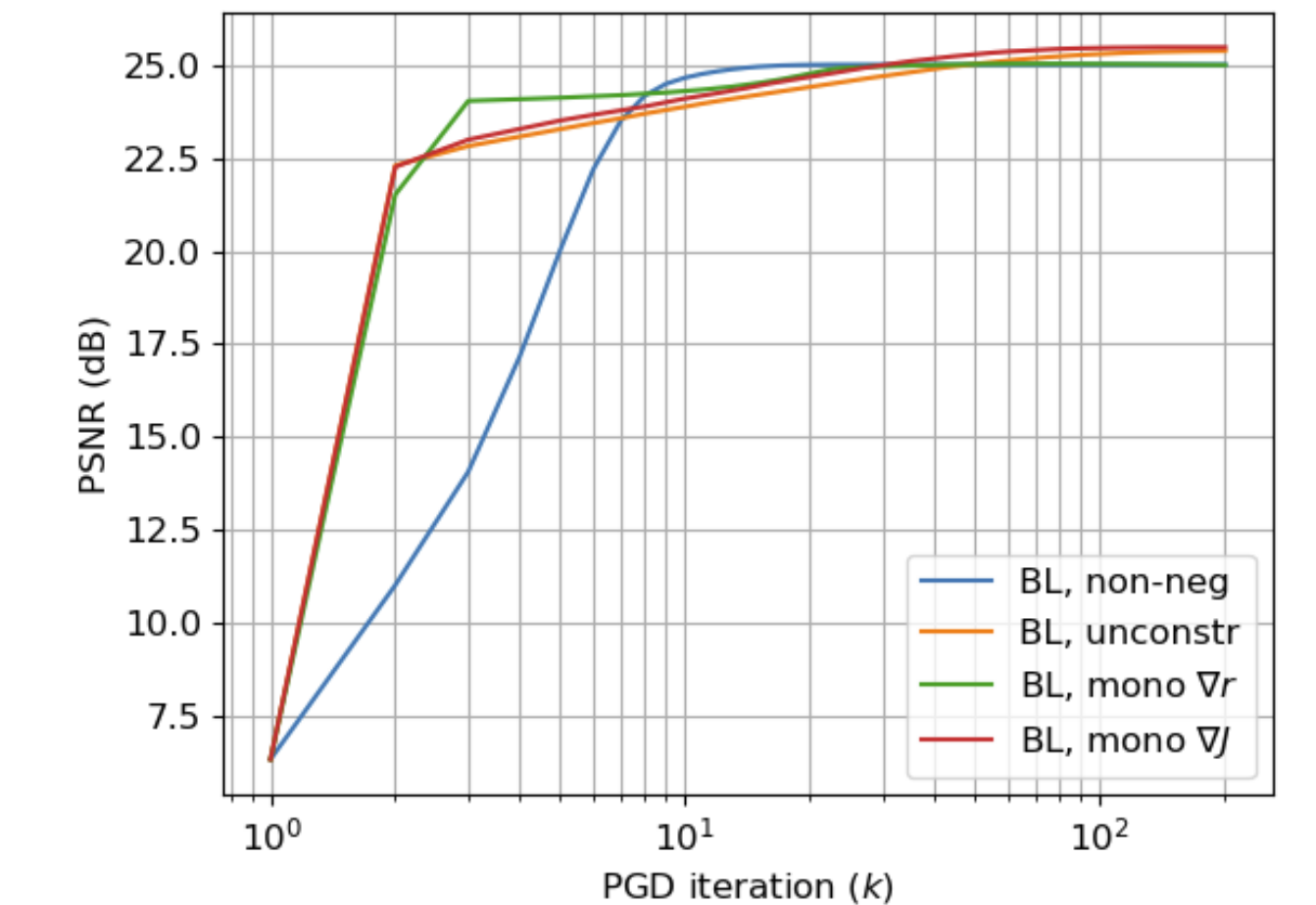
- monotone outperformed structurally convex (more expressive)
- monotone outperformed unconstrained (trapped in local minima?)
- monotone + Lipschitz performed best



Network	PSNR(dB)	SSIM
non-neg	24.77	0.681
unconstrained	25.25	0.698
monotone ∇r	25.28	0.699
mono, 20-Lip ∇r	25.34	0.701

Bi-level (BL) training

- outperforms GS
- monotone- ∇r tied structurally convex
- monotone- ∇J outperformed unconstrained



Network	PSNR(dB)	SSIM
non-neg	25.01	0.683
unconstrained	25.39	0.723
monotone ∇r	25.01	0.692
monotone ∇J	25.48	0.726

Future Work

- Enforce monotone constraint in the vicinity of ground truth data to increase expressivity.
- Study our methods' performance on other regularizer architectures like [9].

References

- S. Roth and M. J. Black, "Field of experts," *Int. J. Comput. Vision*, vol. 82, pp. 205–229, Apr. 2009.
- Y. Romano, M. Elad, and P. Milanfar, "The little engine that could: Regularization by denoising (RED)," *SIAM J. Imag. Sci.*, vol. 10, no. 4, pp. 1804–1844, 2017.
- T. Salimans and J. Ho, "Should EBM model the energy or the score?," in *ICLR Workshop on Energy-based Models*, 2021.
- H. Li, J. Schwab, S. Antholzer, and M. Haltmeier, "NETT: Solving inverse problems with deep neural networks," *Inverse Problems*, vol. 36, no. 6, p. 065005, 2020.
- E. Kobler, A. Effland, K. Kunisch, and T. Pock, "Total deep variation for linear inverse problems," in *Proc. IEEE Conf. Comp. Vision Pattern Recog.*, pp. 7549–7558, 2020.
- R. Cohen, Y. Blau, D. Freedman, and E. Rivlin, "It has potential: Gradient-driven denoisers for convergent solutions to inverse problems," in *Proc. Neural Inf. Process. Syst. Conf.*, vol. 34, pp. 18152–18164, 2021.
- M. Zach, F. Knoll, and T. Pock, "Stable deep MRI reconstruction using generative priors," *arXiv:2210.13834*, 2022.
- S. Mukherjee, S. Dittmer, Z. Shumaylov, S. Lunz, O. Öktem, and C.-B. Schönlieb, "Learned convex regularizers for inverse problems," *arXiv:2008.02839*, 2020.
- A. Goujon, S. Neumayer, P. Bohra, S. Ducotterd, and M. Unser, "A neural-network-based convex regularizer for image reconstruction," *arXiv:2211.12461*, p. 4, 2022.
- S. Hurault, A. Leclaire, and N. Papadakis, "Proximal denoiser for convergent plug-and-play optimization with nonconvex regularization," in *Proc. Int. Conf. Mach. Learn.*, 2022.
- K. G. Samul and M. F. Tappen, "Learning optimized MAP estimates in continuously-valued MRF models," in *Proc. IEEE Conf. Comp. Vision Pattern Recog.*, pp. 477–484, 2009.
- J.-C. Pesquet, A. Repetti, M. Terris, and Y. Wiaux, "Learning maximally monotone operators for image recovery," *SIAM J. Imag. Sci.*, vol. 14, no. 3, pp. 1206–1237, 2021.
- S. V. Venkatakrishnan, C. A. Bouman, and B. Wohlberg, "Plug-and-play priors for model based reconstruction," in *Proc. IEEE Global Conf. Signal Inf. Process.*, pp. 945–948, 2013.
- I. Selesnick, A. Lanza, S. Morigi, and F. Sgallari, "Non-convex total variation regularization for convex denoising of signals," *J. Math. Imaging Vis.*, vol. 62, no. 6, pp. 825–841, 2020.
- S. Elfving, E. Uchibe, and K. Doya, "Sigmoid-weighted linear units for neural network function approximation in reinforcement learning," *Neural Networks*, vol. 107, 2018.