



This paper proposes a deep unfolding model-driven approach for hypersharpening, which involves merging a hyperspectral (HS) image and a panchromatic (PAN) image. The key contributions of this work can be summarized as follows: firstly, a novel model-based method is introduced for hypersharpening by incorporating a constraint for injecting high-frequency details from the panchromatic image into the fused image. Secondly, the model is formulated as an energy functional optimization problem, solved using a forward-backward splitting algorithm. The solution is further unfolded into a deep learning framework.

### Mathematical Model

The observation models [5, 4] relating  ${f H}$  and  ${f P}$  with  ${f U}$  are generally given by

$$egin{aligned} \mathbf{H} &= \mathbf{D}\mathbf{B}\mathbf{U} + oldsymbol{\eta}_h \ \mathbf{P} &= \mathbf{U}\mathbf{S} + oldsymbol{\eta}_p, \end{aligned}$$

where 
$$\mathbf{B} \in \mathbb{R}^{N \times N}$$
 is the low-pass filter modeling the point spread function of the HS sensors,  $\mathbf{D}$  pling operator,  $\mathbf{S} \in \mathbb{R}^{C \times 1}$  is the spectral response of the PAN sensor, and  $\boldsymbol{\eta}_h$  and  $\boldsymbol{\eta}_p$  are assume noise.

Then,  $\mathbf{U}$  can be estimated by solving the following minimization problem:

$$\min_{\mathbf{U}} \frac{1}{2} \|\mathbf{D}\mathbf{B}\mathbf{U} - \mathbf{H}\|^2 + \frac{\gamma}{2} \|\mathbf{U}\mathbf{S} - \mathbf{P}\|^2 + \mu R(\mathbf{U}),$$

On one hand, the high frequencies of the fused image can be estimated as  $\mathbf{U} - \mathbf{ ilde{H}}$ , where  $\mathbf{ ilde{H}} \in \mathbb{R}^{N imes C}$  is the result of upscaling **H** by bicubic interpolation. On the other hand, the high frequencies of the scene are given by  $\mathbf{P} - \tilde{\mathbf{P}}$ , where  $\tilde{\mathbf{P}} \in \mathbb{R}^{N imes 1}$  contains the low frequencies of the PAN data.

$$U_{ij} - \tilde{H}_{ij} = \frac{\tilde{H}_{ij}}{\tilde{P}_i} (P_i - \tilde{P}_i),$$

We can assume that U can be linearly represented by P and an unknown matrix  $\mathbf{V} \in \mathbb{R}^{N \times (r-1)}$ , where  $r = \text{rank}(\mathbf{U}) > 1$ , i.e.,  $\mathbf{U} = \mathbf{P}\mathbf{X} + \mathbf{V}\mathbf{Y}$ 

Putting it all together, the proposed fusion model is

$$\min_{\mathbf{V}} \mu R(\mathbf{V}) + F(\mathbf{V}) = \min_{\mathbf{V}} \mu R(\mathbf{V}) + \frac{1}{2} \|\mathbf{D}\mathbf{B}(\mathbf{P}\mathbf{X} + \mathbf{V}\mathbf{Y}) - \mathbf{H}\|^{2} + \frac{\lambda}{2} \|\mathbf{\tilde{P}} \circ (\mathbf{P}\mathbf{X} + \mathbf{V}\mathbf{Y}) - \mathbf{P} \circ \mathbf{\tilde{H}}\|^{2},$$

### **Minimitzation Method**

To solve (5) we use the forward-backward splitting method [2, 1] to compute the solution. The basic idea is to combine an explicit step of descent in the smooth function F with an implicit step of descent in  $\mu R$ :

$$\mathbf{V}^{k+1} = \operatorname{prox}_{\tau\mu} \left( \mathbf{V}^k - \tau \nabla F(\mathbf{V}^k) \right),$$

where

$$\nabla F(\mathbf{V}) = \mathbf{B}^{\top} \mathbf{D}^{\top} \left( \mathbf{D} \mathbf{B} \left( \mathbf{P} \mathbf{X} + \mathbf{V} \mathbf{Y} \right) - \mathbf{H} \right) \mathbf{Y}^{\top} + \lambda \left[ \mathbf{\tilde{P}} \circ \left( \mathbf{\tilde{P}} \circ \left( \mathbf{P} \mathbf{X} + \mathbf{V} \mathbf{Y} \right) - \mathbf{P} \circ \mathbf{\tilde{H}} \right) \right] \mathbf{Y}^{\top}.$$

Algorithmic Steps	Deep Learning
$\mathbf{U}^{(k)} = \mathbf{P}\mathbf{X} + \mathbf{V}^{(k)}\mathbf{Y}$	$\mathcal{U} = \mathbf{Conv}_{1 \to C}(\mathcal{P}) + \mathbf{Conv}_{(r-1) \to C}(\mathcal{V})$
$\begin{aligned} \mathbf{F}^{(k)} &= \mathbf{D}\mathbf{B}\mathbf{U}^{(k)} - \mathbf{H} \\ \mathbf{J}^{(k)} &= \mathbf{B}^{\top}\mathbf{D}^{\top}\mathbf{F}^{(k)}\mathbf{Y}^{\top} \end{aligned}$	$ \begin{array}{l} \mathcal{F} = \mathbf{dSamp}_{N \rightarrow n}(\mathcal{U}) - \mathcal{H} \\ \mathcal{J} = \mathbf{Conv}_{C \rightarrow (r-1)}(\mathbf{uSamp}_{n \rightarrow N}(\mathcal{F})) \end{array} \end{array} $
$\begin{split} \mathbf{L}^{(k)} &= \tilde{\mathbf{P}}^{\text{col}} \circ \mathbf{U}^{(k)} - \mathbf{P}^{\text{col}} \circ \tilde{\mathbf{H}} \\ \mathbf{T}^{(k)} &= \lambda (\tilde{\mathbf{P}}^{\text{col}} \circ \mathbf{L}^{(k)}) \mathbf{Y}^{\top} \end{split}$	$ \begin{aligned} \mathcal{L} &= \mathbf{pMult}(\tilde{\mathcal{P}}^{\mathrm{col}},\mathcal{U}) - \mathbf{pMult}(\mathcal{P}^{\mathrm{col}},\tilde{\mathcal{H}}) \\ \mathcal{T} &= \mathbf{Conv}_{C \to (r-1)}(\mathbf{pMult}(\tilde{\mathcal{P}}^{\mathrm{col}},\mathcal{L})) \end{aligned} $
$\mathbf{V}^{(k+1)} = \operatorname{Prox}_{\tau\mu}(\mathbf{V}^{(k)} - \tau(\mathbf{J}^{(k)} + \mathbf{T}^{(k)}))$	$\mathcal{V} = \mathbf{ProxNet}(\mathcal{V} - \tau(\mathcal{J} + \mathcal{T}))$

Figure 1. Relationship between the steps of the optimization algorithm and the modules of the deep unfolded network. The operator **ProxNet** stands for the proximal operator and it is replaced by a ResNet [3] as suggested in [7] and **uSamp**<sub>nin</sub> and **dSamp**<sub>nin</sub> are the learned upsampling and downsampling operators.

# Deep unfolding for hypersharpening using a high-frequency injection module

Jamila Mifdal<sup>1</sup>, Marc Tomás-Cruz<sup>2</sup>, Alessandro Sebastianelli<sup>1</sup>, Bartomeu Coll<sup>2</sup>, Joan Duran<sup>2</sup>

<sup>1</sup> European Space Agency (ESA) <sup>2</sup> Universitat de les Illes Balears (UIB) & Institute of Applied Computing Community Code (IAC3)

# **Deep Learning Model**

The right side of Figure 1 shows the corresponding operations in the DL framework of the hypersharpening algorithm. The four blocs highlighted in Figure 1 are the main components of the complete network illustrated in Figure 2 where we could see three main stages. Each stage is composed of the blocks detailed in Figure 1 and at each epoch all three stages are executed, then, the estimated result is fed to the following loss function.





Figure 2. The unfolded network of the hypersharpening algorihtm using a high-frequency injection module. The network is composed of three stages, each one of these stages follows the unfolding steps detailed in Figure 1.

# Implementation details for PRISMA dataset

Our model is trained in a PyTorch framework, using an Nvidia A100 GPU, during 4500 epochs for the PRISMA dataset. We use an Adam optimizer with a learning rate of  $10^{-3}$  and a batch size of 8 images. The trade-off parameters  $\alpha$  and  $\beta$  are fixed at  $10^{-3}$ .

The PAN data contains one single band at a spatial resolution of 5m. We selected and downloaded 20 large-scale scenes of PRISMA images throughout the PRISMA mission's portal (https://prisma.asi.it). The downloaded HS images have an original size of  $1000 \times 1000 \times 240$ , each one of the scenes was cropped into non-overlapping tiles of  $128 \times 128 \times 240$ . Given that the SWIR bands are not covered by the spectral response of the PAN sensor, only the first 66 bands were considered which resulted in tiles of  $128 \times 128 \times 66$ , from each tile a new HS and PAN images were generated following the Wald protocol [6] and using the spectral and spatial responses provided by PRISMA mission engineers. The chosen downsampling factor for the PRISMA dataset is 12, thus, from each tile of  $128 \times 128 \times 66$ , an HS image of  $11 \times 11 \times 66$  and PAN image of the size  $128 \times 128$  were considered for the fusion process which is called hyper-sharpening.

# **Quantitative Results**

	$ERGAS \downarrow PSNR \uparrow SSIM \uparrow DD \downarrow SAM \downarrow$
PCA	376.06 15.70 0.3990 0.1333 35.41
Brovey	92.86 27.68 0.9180 0.0309 4.69
Bicubic	225.81 23.40 0.8303 0.0420 4.64
GS	92.14 27.75 0.9181 0.0308 4.78
GSA	186.01 23.98 0.8706 0.0399 4.67
IHS	101.00 26.78 0.8882 0.0346 7.20
SFIM	208.91 23.95 0.8801 0.0392 4.67
DICNN	41.44 33.38 0.9520 0.0145 3.70
MSDCNN	43.10 33.07 0.9496 0.0153 4.06
GPPNN	253.18 20.99 0.8453 0.0700 7.92
MHFnet	45.29 32.69 0.9402 0.0157 4.13
Ours	15.31 42.17 0.9900 0.0078 1.59

 $\mathbf{D} \in \mathbb{R}^{n \times N}$  is the *l*-fold downsamed to be additive, white Gaussian

### (2)

The best results are in bold and the second best ones are underlined.

Table 1. Average of the quality measures over all images of the PRISMA dataset. The methods are divided in classical, pure DL and unfolding methods.

The visualization results of the model on the Prisma dataset are depicted in Figure 3. These visual outputs are complemented by quantitative findings presented in Table 1. The qualitative analysis further supports the quantitative results, establishing a comprehensive understanding of the model's performance.



(8)



Figure 3. Visual comparison of the fusion approaches on an image of the PRISMA dataset. We display the 35th, 45th and the 57th bands in place of the RGB channels. The proposed deep unfolding network successfully combines the geometry of the PAN image with the spectral information of the HS data, while all other results are affected by blur, color artifacts and spatial distortions. Our method is also able to recover both large structures, such as the circular ground contours, and the finest ones, such as roads and small building structures.

This work is part of the MaLiSat project TED2021-132644B-IOO funded by MCIN/AEI/10.13039/501100011033 and by the European Union "NextGenerationEU"/PRTR. The authors were also supported by the Conselleria de Fons Europeus, Universitat i Cultura del Govern de les Illes Balears under grant AP\_2021\_023.

[1] Antonin Chambolle and Thomas Pock. "An introduction to continuous op- [5] Rafael Molina, Aggelos K Katsaggelos, and Javier Mateos. "Bayesian and timization for imaging". In: Acta Numerica 25 (2016), pp. 161–319.

- [2] Patrick L Combettes and Valérie R Wajs. "Signal recovery by proximal tion". In: IEEE Transactions on Image Processing 8.2 (1999), pp. 231–246. forward-backward splitting". In: Multiscale modeling & simulation 4.4 (2005), [6] Lucien Wald, Thierry Ranchin, and Marc Mangolini. "Fusion of satellite images of different spatial resolutions: Assessing the quality of resulting pp. 1168–1200. images". In: Photogrammetric engineering and remote sensing 63.6 (1997), pp. 691–699.
- [3] Kaiming He et al. "Deep residual learning for image recognition". In: Proceedings of the IEEE conference on computer vision and pattern recognition. 2016, pp. 770-778. [7] Qi Xie et al. "Multispectral and hyperspectral image fusion by MS/HS fu-
- sion net". In: Proceedings of the IEEE/CVF Conference on Computer Vision [4] Jamila Mifdal et al. "Variational Fusion of Hyperspectral Data by Non-Local Filtering". In: *Mathematics* 9.11 (2021), p. 1265. and Pattern Recognition. 2019, pp. 1585–1594.



# **Qualitative Results**

HS	GSA	IHS	Brovey	SFIM
1HFnet	DiCNN	GPPNN	MSDCNN	Ours

# Acknowledgements

## References

regularization methods for hyperparameter estimation in image restora-